

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2004-326801

(P2004-326801A)

(43) 公開日 平成16年11月18日(2004. 11. 18)

(51) Int. Cl.<sup>7</sup>

G06F 12/00

G06F 3/06

F I

G06F 12/00

5 1 4 E

G06F 3/06

3 0 1 K

G06F 3/06

3 0 4 F

テーマコード (参考)

5 B 0 6 5

5 B 0 8 2

審査請求 有 請求項の数 55 O L (全 39 頁)

(21) 出願番号 特願2004-129471 (P2004-129471)  
 (22) 出願日 平成16年4月26日 (2004. 4. 26)  
 (31) 優先権主張番号 10/427403  
 (32) 優先日 平成15年4月29日 (2003. 4. 29)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 390009531  
 インターナショナル・ビジネス・マシー  
 ズ・コーポレーション  
 INTERNATIONAL BUSIN  
 ESS MACHINES CORPO  
 RATION  
 アメリカ合衆国10504 ニューヨーク  
 州 アーモンク ニュー オーチャード  
 ロード  
 (74) 代理人 100086243  
 弁理士 坂口 博  
 (74) 代理人 100091568  
 弁理士 市位 嘉宏  
 (74) 代理人 100108501  
 弁理士 上野 剛史

最終頁に続く

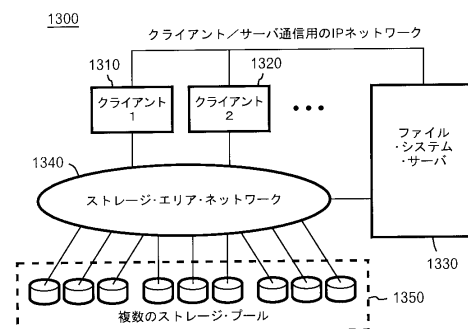
(54) 【発明の名称】 ファイルのコピー・オン・ライトを実施する方法、システム、およびコンピュータ・プログラム

(57) 【要約】 (修正有)

【課題】 コンピューティング環境のさまざまなコピー・オン・ライト実施形態を提示すること。

【解決手段】 あるコピー・オン・ライト実施形態に、修正のために物理記憶装置からコピー・オン・ライトされるファイルのデータのブロックを読み取るのに使用される第1仮想ブロック - 物理ブロック・マッピングを実行するのに読取マッピング・テーブルを使用することと、ファイル・データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック - 物理ブロック・マッピングを実行するのに、異なる書込マッピング・テーブルを使用することとが含まれ、ここで、データのブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される。

【選択図】 図13



**【特許請求の範囲】****【請求項 1】**

コンピューティング環境でコピー・オン・ライトを実施する方法であって、

( i ) ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第 1 仮想ブロック - 物理ブロック・マッピングを実行するのに第 1 マッピング・テーブルを使用することと、

( i i ) 前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第 2 仮想ブロック - 物理ブロック・マッピングを実行するのに第 2 マッピング・テーブルを使用することであって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、前記使用することと

10

を含む方法。

**【請求項 2】**

前記第 1 マッピング・テーブルが、読取マッピング・テーブルを含み、前記第 2 マッピング・テーブルが、書込マッピング・テーブルを含み、前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック - 物理ブロック・マッピングと異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも 1 つの仮想ブロックを含む、請求項 1 に記載の方法。

**【請求項 3】**

前記コピー・オン・ライト実施形態が、さらに、まず修正が部分ブロック書込または全ブロック書込のどちらを含むかを判定することと、部分ブロック書込の場合に、前記使用すること ( i ) および前記使用すること ( i i ) を実行し、そうでない場合に前記使用すること ( i ) を実行せずに前記実行すること ( i i ) を実行することとを含む、請求項 1 に記載の方法。

20

**【請求項 4】**

前記実行すること ( i ) が、データの前記ブロックを物理記憶装置からバッファに読み取ることを含み、前記方法が、さらに、前記使用すること ( i i ) を実行する前に前記バッファ内のデータの前記ブロックを修正することを含む、請求項 1 に記載の方法。

**【請求項 5】**

前記コンピューティング環境が、ファイルシステム・サーバおよび少なくとも 1 つのクライアントを含むクライアント・サーバ環境を含み、前記使用すること ( i ) および前記使用すること ( i i ) が、前記クライアント・サーバ環境の前記少なくとも 1 つのクライアントによって実行される、請求項 1 に記載の方法。

30

**【請求項 6】**

前記少なくとも 1 つのクライアントによって実行される前記使用すること ( i ) および前記使用すること ( i i ) が、前記ファイルのデータの前記ブロックのコピー・オン・ライトを実行する時に、前記第 1 マッピング・テーブルおよび前記第 2 マッピング・テーブルの少なくとも 1 つを入手するために前記ファイルシステム・サーバへの少なくとも 1 つの呼出しを行うことをさらに含む、請求項 5 に記載の方法。

**【請求項 7】**

データの前記修正されたブロックを物理記憶装置に書き込んだ後に、前記第 1 マッピング・テーブルを更新することをさらに含み、前記更新することが、前記第 1 マッピング・テーブルの少なくとも 1 つの仮想ブロック - 物理ブロック変換を、前記第 2 マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と一致するように修正することを含む、請求項 5 に記載の方法。

40

**【請求項 8】**

前記コンピューティング環境が、コンピューティング・ユニットおよび外部記憶ユニットを含み、前記外部記憶ユニットが、前記物理記憶装置を含み、前記使用すること ( i ) および前記使用すること ( i i ) が、前記コンピューティング・ユニットによって実行される、請求項 1 に記載の方法。

**【請求項 9】**

50

クライアント・サーバ・コンピューティング環境でコピー・オン・ライトを容易にする方法であって、

前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持すること

を含み、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第1仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能である

10

方法。

【請求項10】

前記ファイルに関する前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも1つの仮想ブロックを含む、請求項9に記載の方法。

【請求項11】

前記ファイルのデータのブロックのコピー・オン・ライトが実行された後に前記読取マッピング・テーブルを更新することをさらに含み、前記更新することが、前記読取マッピング・テーブルの少なくとも1つの仮想ブロック - 物理ブロック変換を、前記書込マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と一致するように修正することを含む、請求項9に記載の方法。

20

【請求項12】

複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施する方法であって、

前記クライアント・サーバ環境の複数のクライアントを使用してファイルのコピー・オン・ライトを実行することであって、

(i) 前記複数のクライアントの第1クライアントによって、前記ファイルのデータの少なくとも1つのブロックのコピー・オン・ライトを実行することと、

30

(ii) 前記複数のクライアントの第2クライアントによって、前記ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトを実行することと

を含む、実行すること

を含む方法。

【請求項13】

前記実行すること(i)が、前記第1クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行することを含み、前記実行すること(ii)が、前記第2クライアントによって、単一の書込動作を使用してデータの前記少なくとも1つの他のブロックのコピー・オン・ライトを実行することを含む、請求項12に記載の方法。

40

【請求項14】

前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記ファイルシステム・サーバが、前記少なくとも1つの共用記憶装置に保管された前記ファイルに関するマッピング・テーブルを維持し、前記実行すること(i)が、前記第1クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを入手することと、前記ファイルのデータの前記少なくとも1つのブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用することとを含み、前記実行すること(ii)が、前記第2クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する前記読取マ

50

ッピング・テーブルおよび前記書込マッピング・テーブルを入手することと、前記ファイルのデータの前記少なくとも1つの他のブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用することを含む、請求項13に記載の方法。

【請求項15】

前記第1クライアントによって、前記第1クライアントが前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行したことを前記ファイルシステム・サーバに知らせることと、それに応答して、前記ファイルシステム・サーバによって維持される前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも1つを更新することとをさらに含む、請求項14に記載の方法 10

【請求項16】

前記ファイルシステム・サーバが、前記第1クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも1つのブロックおよび前記第2クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも1つの他のブロックに対するコピー・オン・ライトを前記複数のクライアントのすべてのクライアントが行えなくする、請求項14に記載の方法。

【請求項17】

前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記方法が、さらに、前記ファイルの前記コピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御することを含み、前記開始することが、前記実行すること(i)および前記実行すること(ii)によって使用される前記ファイルに関する書込マッピング・テーブルを更新することを含む、請求項12に記載の方法。 20

【請求項18】

クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にする方法であって、

前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御することであって、前記制御することが、前記クライアント・サーバ環境の第1クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにすることと、前記クライアント・サーバ環境の第2クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにすることとを含み、前記ファイルシステム・サーバが、前記ファイルの分散コピー・オン・ライトの実行を制御すること 30

を含む方法。

【請求項19】

前記制御することが、前記ファイルのコピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御することを含み、前記開始することが、前記コピー・オン・ライトを実行するのに使用される前記ファイルに関する少なくとも1つのマッピング・テーブルを更新することを含む、請求項18に記載の方法。 40

【請求項20】

前記ファイルシステム・サーバで、前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持することとをさらに含み、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルが、前記コピー・オン・ライトの実行に使用される、請求項18に記載の方法。

【請求項21】

前記コピー・オン・ライトの実行の後に、前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも1つを更新することとをさらに含む、請求項20に記載の方法。

【請求項22】

前記制御することが、前記ファイルシステム・サーバによって、前記ファイルの前記コピー・オン・ライトの一部としての、前記第1クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記部分に対する追加更新および前記第2クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記異なる部分に対する追加更新を防ぐことをさらに含む、請求項18に記載の方法。

【請求項23】

コンピューティング環境でコピー・オン・ライトを実施するシステムであって、

(i) ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第1仮想ブロック・物理ブロック・マッピングを実行するのに第1マッピング・テーブルを使用する手段と、

(ii) 前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに第2マッピング・テーブルを使用する手段であって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、使用する手段とを含むシステム。

【請求項24】

前記第1マッピング・テーブルが、読取マッピング・テーブルを含み、前記第2マッピング・テーブルが、書込マッピング・テーブルを含み、前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック・物理ブロック・マッピングと異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも1つの仮想ブロックを含む、請求項23に記載のシステム。

【請求項25】

前記コピー・オン・ライト実施形態が、さらに、まず修正が部分ブロック書込または全ブロック書込のどちらを含むかを判定し、部分ブロック書込の場合に、前記使用すること(i)および前記使用すること(ii)を実行し、そうでない場合に前記使用すること(i)を実行せずに前記実行すること(ii)を実行する手段を含む、請求項23に記載のシステム。

【請求項26】

前記実行する手段(i)が、データの前記ブロックを物理記憶装置からバッファに読み取る手段を含み、前記システムが、さらに、前記使用すること(ii)を実行する前に前記バッファ内の前記ブロックを修正する手段を含む、請求項23に記載のシステム。

【請求項27】

前記コンピューティング環境が、ファイルシステム・サーバおよび少なくとも1つのクライアントを含むクライアント・サーバ環境を含み、前記使用する手段(i)および前記使用する手段(ii)が、前記クライアント・サーバ環境の前記少なくとも1つのクライアントによって実行される、請求項23に記載のシステム。

【請求項28】

前記少なくとも1つのクライアントによって実行される前記使用する手段(i)および前記使用する手段(ii)が、前記ファイルのデータの前記ブロックのコピー・オン・ライトを実行する時に、前記第1マッピング・テーブルおよび前記第2マッピング・テーブルの少なくとも1つを入手するために前記ファイルシステム・サーバへの少なくとも1つの呼出しを行う手段をさらに含む、請求項27に記載のシステム。

【請求項29】

データの前記修正されたブロックを物理記憶装置に書き込んだ後に、前記第1マッピング・テーブルを更新する手段をさらに含み、前記更新する手段が、前記第1マッピング・テーブルの少なくとも1つの仮想ブロック・物理ブロック変換を、前記第2マッピング・テーブルの対応する仮想ブロック・物理ブロック変換と一致するように修正する手段を含む、請求項27に記載のシステム。

【請求項30】

前記コンピューティング環境が、コンピューティング・ユニットおよび外部記憶ユニッ

10

20

30

40

50

トを含み、前記外部記憶ユニットが、前記物理記憶装置を含み、前記使用する手段 ( i ) および前記使用する手段 ( i i ) が、前記コンピューティング・ユニットによって実行される、請求項 2 3 に記載のシステム。

【請求項 3 1】

クライアント・サーバ・コンピュータ環境でコピー・オン・ライトを容易にするシステムであって、

前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する手段

を含み、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第 1 仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第 2 仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能であるシステム。 10

【請求項 3 2】

前記ファイルに関する前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも 1 つの仮想ブロックを含む、請求項 3 1 に記載のシステム。 20

【請求項 3 3】

前記ファイルのデータのブロックのコピー・オン・ライトが実行された後に前記読取マッピング・テーブルを更新する手段をさらに含み、前記更新する手段が、前記読取マッピング・テーブルの少なくとも 1 つの仮想ブロック - 物理ブロック変換を、前記書込マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と一致するように修正する手段を含む、請求項 3 1 に記載のシステム。

【請求項 3 4】

複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施するシステムであって、 30

( i ) 前記クライアント・サーバ環境の第 1 クライアントで、コピー・オン・ライトされる前記ファイルのデータの少なくとも 1 つのブロックのコピー・オン・ライトを実行する手段と、

( i i ) 前記クライアント・サーバ環境の第 2 クライアントで、コピー・オン・ライトされる前記ファイルのデータの少なくとも 1 つの他のブロックのコピー・オン・ライトを実行する手段であって、前記ファイルの前記コピー・オン・ライトの異なる部分が、前記クライアント・サーバ環境内の前記複数のクライアントの異なるクライアントによって実行される、手段と

を含むシステム。 40

【請求項 3 5】

前記実行する手段 ( i ) が、前記第 1 クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも 1 つのブロックのコピー・オン・ライトを実行する手段を含み、前記実行する手段 ( i i ) が、前記第 2 クライアントによって、単一の書込動作を使用してデータの前記少なくとも 1 つの他のブロックのコピー・オン・ライトを実行する手段を含む、請求項 3 4 に記載のシステム。

【請求項 3 6】

前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも 1 つの共用記憶装置に関連し、前記ファイルシステム・サーバが、前記少なくとも 1 つの共用記憶装置に保管された前記ファイルに関するマッピング・テーブルを維持し 50

、前記実行する手段 ( i ) が、前記第 1 クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも 1 つのブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する手段を含み、前記実行する手段 ( i i ) が、前記第 2 クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも 1 つの他のブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する手段を含む、請求項 35 に記載のシステム。

10

【請求項 37】

前記第 1 クライアントによって、前記第 1 クライアントが前記ファイルのデータの前記少なくとも 1 つのブロックのコピー・オン・ライトを実行したことを前記ファイルシステム・サーバに知らせ、それに応答して、前記ファイルシステム・サーバによって維持される前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも 1 つを更新する手段をさらに含む、請求項 36 に記載のシステム。

【請求項 38】

前記ファイルシステム・サーバが、前記第 1 クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも 1 つのブロックおよび前記第 2 クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも 1 つの他のブロックに対するコピー・オン・ライトを前記複数のクライアントのいずれもが行えなくする、請求項 36 に記載のシステム。

20

【請求項 39】

前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも 1 つの共用記憶装置に関連し、前記システムが、さらに、前記ファイルの前記コピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御する手段を含み、前記開始する手段が、前記実行する手段 ( i ) および前記実行する手段 ( i i ) によって使用される前記ファイルに関する書込マッピング・テーブルを更新する手段を含む、請求項 34 に記載のシステム。

【請求項 40】

クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にするシステムであって、

30

前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御する手段であって、前記制御する手段が、前記クライアント・サーバ環境の第 1 クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにし、前記クライアント・サーバ環境の第 2 クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにする手段を含み、前記ファイルシステム・サーバが、前記ファイルの分散コピー・オン・ライトの実行を制御し、容易にする、制御する手段

40

を含むシステム。

【請求項 41】

前記制御する手段が、前記ファイルのコピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御する手段を含み、前記開始することが、前記コピー・オン・ライトを実行するのに使用される前記ファイルに関する少なくとも 1 つのマッピング・テーブルを更新することを含む、請求項 40 に記載のシステム。

【請求項 42】

前記ファイルシステム・サーバで、前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する手段をさらに含み、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルが、前記コピー・オン・ライトの実行に使用される、請求項 40 に記載のシステム。

50

## 【請求項 4 3】

前記コピー・オン・ライトの実行の後に、前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも1つを更新する手段をさらに含む、請求項 4 2 に記載のシステム。

## 【請求項 4 4】

前記制御する手段が、前記ファイルシステム・サーバによって、前記ファイルの前記コピー・オン・ライトの一部としての、前記第1クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記部分に対する追加更新および前記第2クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記異なる部分に対する追加更新を防ぐ手段をさらに含む、請求項 4 0 に記載のシステム。

10

## 【請求項 4 5】

製造品であって、

コンピューティング環境でコピー・オン・ライトを実施するコンピュータ可読プログラム・コード論理を有する少なくとも1つのコンピュータ使用可能媒体

を含み、前記コンピュータ可読プログラム・コード論理が、

( i ) ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第1仮想ブロック - 物理ブロック・マッピングを実行するのに第1マッピング・テーブルを使用する論理と、

( i i ) 前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック - 物理ブロック・マッピングを実行するのに第2マッピング・テーブルを使用する論理であって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、使用する論理と

20

を含む、製造品。

## 【請求項 4 6】

前記コンピューティング環境が、ファイルシステム・サーバおよび少なくとも1つのクライアントを含むクライアント・サーバ環境を含み、前記使用する論理 ( i ) および前記使用する論理 ( i i ) が、前記クライアント・サーバ環境の前記少なくとも1つのクライアントによって実行される、請求項 4 5 に記載の製造品。

## 【請求項 4 7】

前記少なくとも1つのクライアントによって実行される前記使用する論理 ( i ) および前記使用する論理 ( i i ) が、前記ファイルのデータの前記ブロックのコピー・オン・ライトを実行する時に、前記第1マッピング・テーブルおよび前記第2マッピング・テーブルの少なくとも1つを入手するために前記ファイルシステム・サーバへの少なくとも1つの呼出しを行うことをさらに含む、請求項 4 5 に記載の製造品。

30

## 【請求項 4 8】

前記コンピューティング環境が、コンピューティング・ユニットおよび外部記憶ユニットを含み、前記外部記憶ユニットが、前記物理記憶装置を含み、前記使用する論理 ( i ) および前記使用する論理 ( i i ) が、前記コンピューティング・ユニットによって実行される、請求項 4 5 に記載の製造品。

## 【請求項 4 9】

40

製造品であって、

クライアント・サーバ・コンピューティング環境でコピー・オン・ライトを容易にするコンピュータ可読プログラム・コード論理を有する少なくとも1つのコンピュータ使用可能媒体

を含み、前記コンピュータ可読プログラム・コード論理が、

前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する論理であって、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第1仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファ

50



イルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能である、論理

を含む、製造品。

【請求項50】

製造品であって、

複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施するコンピュータ可読プログラム・コード手段を有する少なくとも1つのコンピュータ使用可能媒体

10

を含み、前記コンピュータ可読プログラム・コード手段が、

前記クライアント・サーバ環境の複数のクライアントを使用してファイルのコピー・オン・ライトを実行する論理

を含み、前記実行する論理が、

(i) 前記複数のクライアントの第1クライアントによって、前記ファイルのデータの少なくとも1つのブロックのコピー・オン・ライトを実行する論理と、

(ii) 前記複数のクライアントの第2クライアントによって、前記ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトを実行する論理と

を含む、製造品。

【請求項51】

20

前記論理(i)が、前記第1クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行する論理を含み、前記論理(ii)が、前記第2クライアントによって、単一の書込動作を使用してデータの前記少なくとも1つの他のブロックのコピー・オン・ライトを実行する論理を含む、請求項50に記載の製造品。

【請求項52】

前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記ファイルシステム・サーバが、前記少なくとも1つの共用記憶装置に保管された前記ファイルに関するマッピング・テーブルを維持し、前記論理(i)が、前記第1クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも1つのブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する論理を含み、前記論理(ii)が、前記第2クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも1つの他のブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する論理を含む、請求項51に記載の製造品。

30

【請求項53】

40

製造品であって、

クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にするコンピュータ可読プログラム・コード論理を有する少なくとも1つのコンピュータ使用可能媒体

を含み、前記コンピュータ可読プログラム・コード論理が、

前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御する論理であって、前記制御する論理が、前記クライアント・サーバ環境の第1クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにする論理と、前記クライアント・サーバ環境の第2クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにする論理とを含み、前記ファイルシステム・サーバが、前記ファイルの分散コ

50

ピー・オン・ライトの実行を制御し、容易にする、制御する論理を含む、製造品。

【請求項 5 4】

前記制御する論理が、前記ファイルのコピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御する論理を含み、前記開始することが、前記コピー・オン・ライトを実行するのに使用される前記ファイルに関する少なくとも 1 つのマッピング・テーブルを更新することを含む、請求項 5 3 に記載の製造品。

【請求項 5 5】

前記ファイルシステム・サーバで、前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する論理をさらに含み、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルが、前記コピー・オン・ライトの実行に使用される、請求項 5 3 に記載の方法。

10

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明は、全般的にはコンピューティング環境内のファイルシステム・データ管理に関し、具体的には、さまざまなコンピューティング環境内のファイルシステム・データ・ファイルのコピー・オン・ライト (copy-on-write) を実施する技法に関する。

【背景技術】

【0 0 0 2】

汎用コンピュータおよびデータ処理システムを含む多数のタイプのコンピューティング環境で、「仮想記憶」方式を使用して編成された記憶装置が使用される。一般的な仮想記憶を用いると、コンピューティング環境で実行中のアプリケーションまたはプロセスあるいはその両方が、自由に使える無制限の量のメモリを有するかのようには振る舞えるようになる。実際には、特定のアプリケーションまたはプロセスから使用可能な記憶装置の量は、コンピューティング環境内の記憶装置の量によって制限され、さらに、その記憶装置を共用する並行に実行中のプログラムの数によって制限される。さらに、仮想記憶方式によって、メモリの実際の物理アドレスがアプリケーション・プログラムから隠蔽される。アプリケーション・プログラムは、論理アドレスを使用してそのプログラムのメモリ空間にアクセスし、この論理アドレスが、コンピューティング環境によって物理アドレスに変換される。

20

30

【0 0 0 3】

仮想記憶システムによって、記憶装置が「ブロック」(または「ページ」)と称する単位で編成される。これらのブロックは、高速な主記憶と、1つまたは複数のより大きいが通常は低速の二次、三次などの記憶ユニットの間で移動される。ブロックの移動(しばしばスワップと称する)は、コンピューティング環境で実行されるアプリケーションまたはプロセスには透過的であり、アプリケーションまたはプロセスは、それぞれが無制限の量のメモリを有するかのようには振る舞うことができる。

【0 0 0 4】

ある種の従来のシステムは、時々、メモリの諸部分をコピーすることを必要とする。このコピーは、ユーザが開始するか、オペレーティング・システムによって開始されるのいずれかとすることができる。従来のシステムでは、しばしば、「フラッシュ・コピー (flash

40

copy)」の「遅延 (lazy)」コピー方法が使用され、この場合に、コピーされる記憶装置が、読取専用の状況を割り当てられるが、実際のコピーは後に延期される。オリジナルまたはコピーのいずれかへの書込の試みが行われる場合に、メモリが、その時にコピーされ、オリジナルとコピーの両方に、読取・書込の入出力状況が与えられる。この形で、コピーが即座に行われるように見えるが、実際のコピーは、可能な最後の時まで延期される。書込が実行されない場合には、コピーは行われない。この理由から、この方法が、「コピー・オン・ライト」または「仮想コピー (virtualcopy)」と呼ばれる。

50

## 【 0 0 0 5 】

一般に、コピー・オン・ライト動作は、単一の書込が2つの書込動作をもたらすので、計算的に高価である。すなわち、既存のデータ・ブロックを、古い物理ブロックから新しい物理ブロックにコピーする必要がある、その後、実際の更新／書込動作が、新しい物理ブロックに対して実行される。

【非特許文献1】ランダル・クライトン・バーンズ (Randal ChiltonBurns) 著、「ストレージ・エリア・ネットワーク用の分散ファイル・システムでのデータ管理 (DataManagement In A DistributedFile System For StorageArea Networks)」、カリフォルニア大学 (University of California)、米国サンタクルーズ (SantaCruz)、2000年3月

## 【 発明の開示 】

10

## 【 発明が解決しようとする課題 】

## 【 0 0 0 6 】

この計算オーバヘッドに鑑みて、当技術分野に、二重書込要件を部分的に除去する新規のコピー・オン・ライト実施形態の必要がある。

## 【 課題を解決するための手段 】

## 【 0 0 0 7 】

一態様では、コンピューティング環境でコピー・オン・ライトを実施する方法を介して、従来技術の短所が克服され、追加の利益がもたらされる。この方法には、ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第1仮想ブロック・物理ブロック・マッピングを実行するのに第1マッピング・テーブルを使用することと、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに第2マッピング・テーブルを使用することとであって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、使用することとが含まれる。

20

## 【 0 0 0 8 】

もう1つの態様では、クライアント・サーバ・コンピューティング環境でコピー・オン・ライトを容易にする方法が提供される。この方法には、前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持することが含まれ、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第1仮想ブロック・物理ブロック・マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能である。

30

## 【 0 0 0 9 】

もう1つの態様では、複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施する方法が提供される。この方法には、前記クライアント・サーバ環境の複数のクライアントを使用してファイルのコピー・オン・ライトを実行することが含まれる。実行することにより、前記複数のクライアントの第1クライアントによって、コピー・オン・ライトされる前記ファイルのデータの少なくとも1つのブロックのコピー・オン・ライトを実行することと、前記複数のクライアントの第2クライアントによって、コピー・オン・ライトされる前記ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトを実行することとが含まれる。機能強化された態様では、実行することにより、前記第1クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行することと、前記第2クライアントによって、やはり単一の書込動作を使用してデータの前記少なくとも1つの他のブロックのコピー・オン・ライトを実行することとが含まれる。

40

## 【 0 0 1 0 】

50

もう1つの態様では、クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にする方法が提示される。この方法には、前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御することであって、前記制御することが、前記クライアント・サーバ環境の第1クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにすることと、前記クライアント・サーバ環境の第2クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにすることとを含み、前記ファイルシステム・サーバが、前記ファイルの分散コピー・オン・ライトの実行を制御し、容易にする、制御することが含まれる。

#### 【0011】

10

上で要約した方法のさまざまな追加の特徴および機能強化も、上で要約した方法に対応するシステムおよびコンピュータ・プログラム製品と同様に、本明細書に記載され、請求される。

#### 【0012】

さらに、本発明の教示を介して、追加の特徴および長所が実現される。本発明の他の実施形態および態様は、本明細書で詳細に説明され、請求される発明の一部とみなされる。

#### 【0013】

発明とみなされる主題は、本明細書の先頭の請求項で具体的に指摘され、明確に請求される。本発明の前述および他の目的、特徴、および長所は、添付図面と共に解釈される下記の詳細な説明から明白になる。

20

#### 【発明を実施するための最良の形態】

#### 【0014】

##### 概要

本明細書で提示されるのは、一態様で、コンピューティング環境でコピー・オン・ライトを実施する技法である。この技法には、異なる変換すなわち、読取マッピング・テーブルおよび書込マッピング・テーブルを使用して、単一書込動作を使用してファイル内のデータの単位のコピー・オン・ライトを達成することが含まれる。一例として、コピー・オン・ライトは、第1の仮想ブロック・物理ブロック・マッピングを使用して、修正のために物理記憶装置からファイルのデータのブロックを読み取ることと、その後、第2の仮想ブロック・物理ブロック・マッピングを使用してそのデータの修正されたブロックを物理記憶装置に書き込むことによって達成され、第1仮想ブロック・物理ブロック・マッピングおよび第2仮想ブロック・物理ブロック・マッピングに、異なるマッピングが含まれる。

30

#### 【0015】

もう1つの態様で、本明細書に提示されるのは、クライアント・サーバ環境の複数のクライアントにまたがるファイルの分散コピー・オン・ライトを実施する技法である。そのような環境内で、複数のクライアントの第1クライアントによって、ファイルのデータの少なくとも1ブロックのコピー・オン・ライトが実行され、複数のクライアントの第2クライアントによって、ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトが実行される。一実施形態では、クライアントに、異種オペレーティング・システムを含めることができ、コピー・オン・ライトされるファイル内のデータのブロックのコピー・オン・ライトのそれぞれが、単一の書込動作を使用して実行される。また、コピー・オン・ライトを、上で要約したように第1マッピング変換（たとえば、読取マッピング・テーブルを使用する）および第2マッピング変換（たとえば、書込マッピング・テーブルを使用する）を使用して達成することができる。本発明の上記および他の態様を、下で説明し、請求項に列挙する。

40

#### 【0016】

##### 詳細な説明

全体的に100で示される、本発明の態様によるコピー・オン・ライトが組み込まれ、使用されるコンピューティング環境の例を、図1に示す。図からわかるように、コンピュ

50

ーティング環境 100 には、たとえば、少なくとも 1 つの中央処理装置 102、メモリ 104、および 1 つまたは複数の記憶ユニットまたは記憶装置 106 が含まれる。

【0017】

既知のように、中央処理装置 102 は、コンピューティング・ユニットの制御センタであり、これによって、命令実行、割込みアクション、タイミング機能、初期プログラム・ロード、および他の機械関連機能のシーケンシング機能および処理機能が提供される。中央処理装置によって、少なくとも 1 つのオペレーティング・システムが実行され、オペレーティング・システムは、既知のように、他のプログラムの実行を制御し、周辺装置との通信を制御し、コンピュータ・リソースの使用を制御することによって、コンピューティング・ユニットの動作を制御するのに使用される。

10

【0018】

中央処理装置 (CPU) 102 は、メモリ 104 に結合され、このメモリ 104 は、直接にアドレス可能であり、メモリ 104 によって、中央処理装置によるデータの高速度処理がもたらされる。メモリ 104 には、下でさらに説明するように CPU 102 によって使用されるバッファまたはキャッシュ領域 103 が含まれる。もう 1 つの実施形態では、バッファ 103 を、CPU 102 内に置くことができる。記憶ユニット 106 は、入出力装置の一例である。本明細書で使用する記憶ユニット 106 は、コンピューティング・ユニットの外部またはコンピューティング環境 100 のコンピューティング・ユニット内とすることができ、記憶ユニット 106 には、たとえば、主記憶、磁気記憶媒体 (たとえば、テープ、ディスク)、および直接アクセス記憶装置などを含めることができる。データを、図示のように、CPU 102、メモリ 104、および記憶ユニット 106 の間で転送することができる。

20

【0019】

一例では、コンピューティング環境 100 が、単一のシステム環境であり、AIX オペレーティング・システムで動作する RS/6000 コンピュータ・システムが含まれる (RS/6000 および AIX は、International Business Machines Corporation 社によって提供される)。しかし、本発明は、そのような環境に制限されない。本発明の機能を、多数の他のタイプのコンピュータ環境および多数のタイプのコンピュータ・システムに組み込み、使用することができる。たとえば、コンピュータ環境 100 に、分散コンピューティング環境を含めることができ、UNIX (R) ベースのオペレーティング・システムが動作する UNIX (R) ワークステーションを含めることができる。他の変形形態も、可能であり、請求される発明の一部とみなされる。

30

【0020】

既知のように、ファイルは、ユーザ / アプリケーション・データの保管に使用することができる、コンピューティング環境のファイルシステム内の名前付きオブジェクトである。このデータには、ファイル名、オフセット、および長さを指定することによってアクセスすることができる。ユーザ・アプリケーションまたはプロセスにとって、ファイル上のデータは、連続的に見えるが、記憶ユニット (ディスクなど) 内では、データ表現が異なる可能性がある。各ファイルシステムによって、仮想 (相対) オフセット・ブロック番号から物理ブロック番号への間のマッピングまたは変換を提供するマッピング・テーブルが維持され、ブロックは、ファイル内のデータのページまたは他の単位とすることができ、この単位のサイズは、ファイルシステムによって指定される。

40

【0021】

図 1 の例では、ファイルシステムが、記憶ユニットを含むと仮定され、記憶ユニットは、やはり、コンピューティング環境の特定のコンピューティング・ユニットに外付けまたは内蔵とすることができ、図 2 に、アプリケーションまたはプロセスのデータ範囲およびオフセットを、所与のファイルの仮想 / 相対ブロック番号に相関させるファイル表現 200 の一例を示す。各ファイルが、それ自体の仮想 / 相対ブロック番号の組を有することに留意されたい。この例では、0 から 4 K バイトまでのファイル・データが、仮想 / 相対

50

ブロック番号 1 にマッピングされ、4 K から 8 K バイトまでのデータが、ブロック番号 2 にマッピングされ、8 K から 12 K バイトまでのデータが、ブロック番号 3 にマッピングされ、12 K から 16 K バイトまでのデータが、ブロック番号 4 にマッピングされる。やはり、これらの数は、例としてのみ提供される。

#### 【0022】

図 3 に、特定のファイルのファイルシステム・マッピング・テーブル 300 の一例を示す。テーブル 300 は、ファイルの仮想 / 相対ブロック番号を記憶ユニットの実際の物理ブロック・アドレスに変換する際に使用される。たとえば、仮想 / 相対ブロック番号 1、2、3、および 4 が、記憶ユニット 400 の物理ブロック・アドレス A、D、G、および L にマッピングされることが示されている（図 4 参照）。

10

#### 【0023】

たとえば、アプリケーションまたはプロセスによって、特定のファイルのオフセット 5000 から 4 バイトのデータの読取が望まれる場合に、図 2 および図 3 のファイル表現 200 およびファイルシステム・マッピング・テーブル 300 を使用すると、実際のデータ読取は、図 5 の物理ブロック D から行われることがわかる。これは、オフセット 5000 から始まる 4 バイトのデータが、仮想 / 相対ブロック番号 2 に含まれ、これが、図 3 に示されているように、記憶ユニットの物理ブロック・アドレス D に変換されるからである。

#### 【0024】

最初に注記したように、フラッシュ・コピー動作によって、記憶装置の空間効率のよいコピーをすばやく行えるようになる。この動作は、高速である必要があるので、物理的なコピーは、当初は、動作の一部として作られない。後に、適用可能なファイル・データを修正する試みによって、コピー・オン・ライト動作がもたらされる。クライアント・サーバ環境では、メタデータ・コピー・オン・ライトが、通常は、ファイルシステム・サーバによって実行され、ファイル・データ・コピー・オン・ライトは、クライアントによって実行される。PageIn スレッドおよび PageOut スレッドを使用して、データのブロックをクライアントのキャッシュに持ち込み、そのデータを更新し、その後、そのデータを記憶ユニットに書き戻すことができる。本明細書で説明するように、PageIn および PageOut について異なる変換が使用される場合に、クライアントは、潜在的なコピー・オン・ライト・データをバッファに読み取り、バッファ内のデータに更新を適用し、修正されたデータを、PageOut スレッドを介して記憶ユニットの新しい位置に書き込むことができる。したがって、2 つのマッピング・テーブルまたは変換を有することによって、一実施形態で、既存の PageIn および PageOut の概念を活用するコピー・オン・ライト技法がもたらされる。

20

30

#### 【0025】

図 6 は、書込動作ならびに本発明の態様によるコピー・オン・ライト動作の一実施形態の流れ図実施形態である。書込動作 600 は、ファイルのデータのブロック全体（またはページ全体）を書き込むかどうかを判定すること 610 によって開始される。データの部分ブロック書込の場合には、ファイル内のデータの適用可能なブロックを、記憶ユニットからローカル・バッファに読み取る 620（PageIn）。バッファ内のデータのブロックの更新 630 がこれに続き、その後、ファイルのデータの修正されたブロックの記憶ユニットへの書込 640 が続く。ファイルのデータのブロック全体を書き込む場合には、この論理では、単純に、記憶ユニット（図 6 の 650）へのデータのブロック全体の書込に進む。既知のように、図 6 の論理は、オペレーティング・システム・カーネル内で実施することができる。図 7 および 8 に、図 6 の論理を使用して書込動作を実施するさらなる処理の例を示す。

40

#### 【0026】

図 7 に、ファイルのデータのブロックを記憶ユニットからバッファへ読み取る処理の一実施形態を示す。アプリケーション入力は、やはり、読み取られるファイル・データのオフセットおよび長さであり 700、これは、仮想ブロック番号を計算する 710 のに使用される（たとえば、図 2 に示されたものなどのファイル表現を使用して）。次に、サブル

50

ーチン呼出しFile・Getを実行して、仮想ブロック番号を、記憶ユニット内のファイル・データの物理ブロック・アドレスにマッピングする720。物理ブロック・アドレスを使用して、データのブロックを、記憶ユニットからローカル・バッファに読み取る730。

【0027】

図8に、ファイルのデータのブロックを記憶ユニットに書き込む処理の一例を示す。図からわかるように、仮想ブロック番号800を、File・Getマッピング処理で使用して、たとえば図3に示されたものなどのファイルシステム・マッピング・テーブルを使用して、仮想-物理マッピングを入手する810。この物理ブロック・アドレスが、修正されたデータ・ブロックをバッファから記憶ユニットに書き込む時820に使用される。一実施形態では、図7のブロック読取処理に、PageInスレッド処理を含めることができ、図8のデータ・ブロック書込プロセスに、PageOutスレッド処理を含めることができる。

10

【0028】

有利なことに、本明細書で開示されるのは、図6の高水準書込論理フローの変更なしでコピー・オン・ライトを達成する技法である。この技法では、コピー・オン・ライトが実行されるデータの特定のファイルに関する、読取マッピング・テーブルおよび書込マッピング・テーブルと称する2組の変換またはマッピング・テーブルが使用される。一実施形態で、コピー・オン・ライトが実行される時に、必ず、これらの2つのマッピング・テーブルが、ファイルシステムによって維持され、クライアント・アプリケーションによってアクセスされる。たとえば、この2つのマッピング・テーブルが、物理アドレス変換境界でファイルシステム・ドライバに同時に提示される。

20

【0029】

たとえば、図9に、読取マッピング・テーブル900および書込マッピング・テーブル910を示す。読取マッピング・テーブル900によって、仮想ブロック番号1、2、3、および4が、それぞれ物理ブロック番号A、D、G、およびLにマッピングされ、書込マッピング・テーブル910によって、仮想ブロック番号1、2、3、および4が、それぞれ物理ブロックW、X、Y、およびZにマッピングされる。読取マッピング・テーブルによって、読取動作に使用される第1の仮想-物理変換が提供され、書込マッピング・テーブルによって、書込動作に使用される第2の仮想-物理変換が提供される。具体的に言うと、一例として、コピー・オン・ライトを、PageInについて読取テーブル変換、PageOutについて書込テーブル変換を使用して実施することができる。

30

【0030】

図10に、図6の論理によるコピー・オン・ライト動作に使用されるデータ・ブロック読取処理の実施形態を示す。図からわかるように、アプリケーションによって、ファイル内のデータ・オフセットおよび長さが指定され1000、これから仮想ブロック番号が計算される1010(図2参照)。次に、ファイルに関する読取マッピング・テーブル(たとえば図9のテーブル900)を使用して、仮想-物理読取マッピングを入手する1020。次に、この物理読取マッピングを使用して、ファイルの少なくとも1つのデータ・ブロックを記憶ユニットからローカル・バッファに読み取る1030。

【0031】

図11に、図6の論理によるコピー・オン・ライト動作に関するデータ・ブロック書込処理の例を示す。仮想ブロック番号1100を使用して、ファイルの書込マッピング・テーブル(たとえば図9のテーブル910)を使用して、仮想-物理「書込」マッピング1110を入手する。物理書込マッピングを使用して、修正されたデータ・ブロックを、ローカル・バッファから対応する記憶ユニット物理ブロックに書き込む1120。上で注記したように、ファイルのデータのブロック全体のコピー・オン・ライトが実行される場合に、修正されたデータ・ブロックが、図11の物理「書込」マッピングを使用して、単純に記憶装置に直接に書き込まれる。その場合には、ファイルの書込マッピング・テーブルだけが使用される。

40

【0032】

当業者は、読取、更新、および書込が記憶ユニット内の同一の物理ブロック・アドレス

50

で行われるようにするために、読取マッピング・テーブルおよび書込マッピング・テーブルを同一にすることによって、図 6、10、および 11 の論理を使用して通常の手続きを実行できることに気付くであろう。しかし、有利なことに、この論理は、コピー・オン・ライトについても同一である。コピー・オン・ライトが、ファイルについて要求される時に、ファイルシステムによって、たとえば新しい物理アドレス・ブロックが割り振られ、そのファイルの対応する書込マッピング・テーブルが更新される。それを行うことによって、書込動作が、異なる物理ブロック・アドレスで行われ、読取動作もそのアドレスで行われ、これは、オリジナルのデータが記憶ユニット内で手付かずのままになることを意味する。書込動作の後に、そのファイルの読取マッピング・テーブルを、特定のコピー・オン・ライト・アプリケーションに応じて更新することができる。

10

#### 【0033】

図 9 のマッピング・テーブルを参照し、図 5 の例を使用すると、コピー・オン・ライト変更が、ファイルのオフセット 5000 の 4 バイトのデータに対して行われる（たとえば、データ内容「1 2」を「2 1」に変更する）場合に、データの更新されたブロックが、記憶ユニットの物理ブロック番号 X（書込マッピング・テーブル 910（図 9）を使用すると、仮想ブロック番号 2 に対応する）に書き込まれる。コピー・オン・ライト動作の後に、対応する読取マッピング・テーブル（図 9 の 900）を、仮想ブロック番号 2 が記憶ユニット内の物理ブロック番号 X に変換されるように更新することができ、これを、図 12 の更新された読取マッピング・テーブル 1200 に示す。図 12 の書込マッピング・テーブル 1210 は、図 9 の書込マッピング・テーブル 910 と同一である。しかし、

20

#### 【0034】

図 13 に、全体的に 1300 で示される、本発明の 1 つまたは複数の態様を組み込み、使用することができるコンピューティング環境のもう 1 つの例を示す。環境 1300 には、クライアント 1 1310 およびクライアント 2 1320 を含む複数のクライアントが含まれ、これらのクライアントは、この例では、クライアント / サーバ通信に関するインターネット・プロトコル・ネットワークによってファイルシステム・サーバ 1330 に接続される。サーバ 1330 は、ストレージ・エリア・ネットワーク 1340 に接続され、ストレージ・エリア・ネットワーク 1340 は、ファイル・データの保管に使用可能な複数のストレージ・プール 1350 を有する。クライアント 1 1310 およびクライアント 2 1320 は、ストレージ・エリア・ネットワーク 1340 にも直接に接続される。

30

#### 【0035】

一例として、コンピューティング環境 1300 は、1 つの位置（たとえばファイルシステム・サーバ）での本明細書に記載のものなどのマッピング・ファイル（たとえば、ファイルの読取マッピング・テーブルおよび書込マッピング・テーブル）の維持、クライアント・アプリケーションがストレージ・エリア・ネットワーク（SAN）を介して記憶ユニット（すなわち複数のストレージ・プール）に直接にアクセスできること、およびクライアント・アプリケーションが複数のストレージ・プール内のファイルの任意のオブジェクトへの読取 / 書込アクセスを有することを含む、ある種の特徴を有すると仮定される。そのような SAN 環境は、ランダル・クライトン・バーンズ（Randal Chilton Burns）著、「ストレージ・エリア・ネットワーク用の分散ファイル・システムでのデータ管理（Data Management

40

50



In A Distributed File System For Storage Area Networks)」、カリフォルニア大学 (University of California)、米国サンタクルーズ (Santa Cruz)、2000年3月による理論を含めて、さまざまな刊行物に詳細に記載されている。

【0036】

本発明のもう1つの態様では、異なるクライアント・アプリケーションが1つのファイルの異なる部分を更新する、分散コピー・オン・ライト機能が提示される。たとえば、クライアント1が、ファイルの仮想ブロック2のコピー・オン・ライトを実行し、その間に、クライアント2が、そのファイルの仮想ブロック4のコピー・オン・ライトを実行する。図14および15に、この例を詳細に示す。

【0037】

図14に、仮想ブロック2のコピー・オン・ライトを実行するクライアント1の一例を示す。クライアント1は、まず、ファイルシステム・サーバから、「my file」というラベルを付けられたファイルのマッピング・テーブルを入手する1400。たとえばオフセット5000で更新する間に、コピー・オン・ライトが、仮想ブロック2に対して実行される1410。クライアント1は、仮想ブロック2のコピー・オン・ライトを実行したことをファイルシステム・サーバに知らせ1420、ファイルシステム・サーバが、それ相応にマッピング・テーブルを更新する1430。一例では、コピー・オン・ライトを、図6および9から12に関して上で説明したように実行することができる。

【0038】

クライアント1は、ロック機構を使用して、ファイル・データの仮想ブロック2のコピー・オン・ライト更新を実行する。ロック機構には、ファイルシステム・オブジェクトごとのロックが含まれ、クライアントは、所与のファイルに対して動作を実行するために、ファイルシステム・サーバからそのロックを獲得する必要がある。したがって、一実施形態では、クライアント1は、このロックを受け取る時に、そのファイルのマッピング・テーブルも受け取り、クライアント1がロックを失う時に、そのファイルに関するクライアント1側のすべてのマッピング・テーブルが無効になる。したがって、クライアント1が次にロックを入手する時に、クライアント1は、既存のマッピング・テーブルを使用することができず、その代わりに、ファイルシステム・サーバから現在のマッピング・テーブルを入手する。

【0039】

さらなる説明として、図13のコンピューティング環境で、分散ロック機構を使用して、異なるクライアントからのファイルシステム・オブジェクトへのアクセスを制御することができる。分散ロックには2つのタイプすなわちセッション・ロックとデータ・ロックがある。これらのロックは、ファイルシステム・オブジェクトごとのロックである。セッション・ロックは、オープン「ファイル記述子」と同等である。このロックがクライアントによって獲得された時に、このクライアントがこのファイルを使用することに関心を持っていることがサーバに知らされる。クライアントは、異なるモード、たとえば読取モード、書込モード、排他モードなどでセッション・ロックを獲得することができる。たとえば、クライアントAが、排他モードでセッション・ロックを保持している時に、もう1つのクライアント、クライアントBが、同一のファイルのオープンを望む時に、クライアントBは、セッション・ロックに関する要求をサーバに送り、サーバは、クライアントAが排他モードのロックを有するのでそのロック要求を拒否し、したがって、クライアントBは、ファイルを操作することができない。第2のタイプの分散ロックが、データ・ロックである。データ・ロックは、物理的な読取/書込を行うのに使用される。クライアントは、ファイルに対する読取/書込を行うために、読取/書込モードのデータ・ロックを有しなければならない。さらなる説明として、セッション・ロックを、2つのクライアントが書込モードで同時に保持することができる。これは、両方のクライアントが、特定のファイルに書き込むアクセス権を有することを意味するが、これらのクライアントは、書込モードでデータ・ロックを入手するまでは書き込むことができない。しかし、所与の時点で、1つのクライアントだけが、「書込モード」のデータ・ロックを得ることができ、その

10

20

30

40

50

結果、1つのクライアントだけが、実際の入出力を実行することができる。両方のクライアントが、能動的にファイルに作用している場合に、これらのクライアントは、書込モードのセッション・ロックを保持しているが、データ・ロックは、書込を完了するために両方のクライアントの間を往復する。

#### 【0040】

ある時点で、クライアント2が、ファイル「my file」の別の部分のコピー・オン・ライトを実行するためにロックを要求する。図15からわかるように、クライアント2は、ファイル「myfile」のマッピング・テーブルをファイルシステム・サーバから入手する1500。オフセット14000（すなわち仮想ブロック番号4）でファイル内容を更新する間に、クライアント2は、ファイルの仮想ブロック番号4に対するコピー・オン・ライトを実行する1510。その後、クライアント2は、このコピー・オン・ライト更新に関してファイルシステム・サーバに知らせ、サーバは、それ相応に「myfile」の読取マッピング・テーブルを更新する1530。図16に、図14および15の「my file」のコピー・オン・ライト更新から生じる更新された読取マッピング・テーブル1600および書込マッピング・テーブル1610（図9のマッピング・テーブルから開始された）を示す。図からわかるように、読取マッピング・テーブルは、仮想ブロック番号2を物理ブロック番号Xにマッピングされ、仮想ブロック番号4を物理ブロック番号Zにマッピングされている。

10

#### 【0041】

要約すると、図13に示されたものなどの環境内でのコピー・オン・ライトの実施形態を、複数のクライアント・アプリケーションにまたがって分散させることができる。したがって、複数のクライアント・アプリケーションを、分散ファイルシステムでのファイルのコピー・オン・ライトの実行に用いることができる。上で説明したものなど、1時に1つのクライアント・アプリケーションだけによって所有される、分散ロック機構を使用することができる。たとえば、クライアント1が、分散ロックを有する時に、クライアント1は、コピー・オン・ライト動作に関する有効な読取マッピング・テーブルおよび書込マッピング・テーブルを有する。クライアント2が、同一のファイルの別の部分のコピー・オン・ライトを行うことを望む場合に、ファイルシステム・サーバは、クライアント1から分散ロックを奪い、クライアント2に与える。代替の例では、分散ロックを、ファイル範囲によって分割することができ、その結果、クライアント1およびクライアント2の両方が、ファイルの異なる部分に対するコピー・オン・ライトを同時に実行できるようになる。

20

30

#### 【0042】

図17に、分散クライアント・サーバ環境でのコピー・オン・ライトに関する「データ保存」アプリケーションの一例を示す。まず、管理者は、ファイル「my file」のポイント・イン・タイム・イメージを作ると決定する1700。管理コマンドを実行して、「my file」を保存し1710、次に、ファイルシステム・サーバが、クライアント・アプリケーションからファイル「myfile」の制御を奪う1720。たとえば、ファイルシステム・サーバは、そのファイルに関するすべての分散ロックをクライアント・アプリケーションから奪う（このステップは、非分散非クライアント/サーバ環境では適用可能でない）。次に、ファイルシステム・サーバが、書込マッピング・テーブルを更新してコピー・オン・ライトを促進し1730、サーバが、めいめいのクライアント・アプリケーションに制御を返す1740。この点から後は、クライアント・アプリケーションによる「myfile」に対するすべての更新が、ファイルのその特定の部分に関してコピー・オン・ライトが実行されることを自動的にもたらすと同時に、古いデータ・パスが保存される（すなわち、読取マッピング・テーブルの物理ブロック番号を、初期ファイル参照を保存するために維持することができる）。

40

#### 【0043】

特定の例

本発明のさまざまな態様の1つの詳細な実施形態を、次に提示する。この詳細な説明で

50

は、2タイプの入出力動作すなわち、バッファ入出力および直接入出力が可能である。

#### 【0044】

バッファ付き入出力

バッファ付き入出力は、バッファ・キャッシュを介して実行される入出力を意味する。この場合に、読取／書込は、まずキャッシュに頼り、その後、このキャッシュデータが、記憶ユニット（たとえばディスク）にハードニング（harden）される。

#### 【0045】

既存ファイルに対する更新は、データをキャッシュに読み取ることによって行われる。すべての変更をキャッシング／バッファリングされたページに適用することと、そのデータをディスクに書き戻すこと。この事実を考慮に入れると、本明細書で開示されるコピー・オン・ライトは、追加コストなしで達成される。たとえば、分散ファイルシステムにアクセスするクライアントに、少なくとも2つの構成要素すなわち、（i）オペレーティング・システム固有であり、アプリケーションからの要求を処理し、バッファ・メモリ・システムおよびディスクなどの記憶装置と通信するインストール可能ファイル・システム（IFS）、および（ii）すべてのオペレーティング・システムに共通し、ロック管理を処理し、サーバと通信するクライアント状態マネージャ（CSM）が含まれる。フラッシュ・コピーの展望から、IFSは、下記の動作についてCSMに接触する。

- ・読取 CSM API、`csmTranslateBlocks()`を使用して、仮想 - 物理変換を得る。
- ・書込 下記の2ステップで行われる

1. IFSは、ページ・キャッシュへの書込を受け入れる前に、バッキング・ブロックを有することを確かめる必要がある。したがって、IFSによって、API、`csmAttachBlocks()`を介してCSMが呼び出される。

2. `csmAttachBlocks`の成功の際に、IFSによって書込を完了できるようになる。

- ・DirectIO：この場合に、書込は直接にディスクに対して行われる。
- ・BufferedIO：この場合に、IFSは、ページをキャッシュに取り込み、それを修正し、その後、ディスクに書き込む。

- ・切捨 CSM API、`csmDetachBlocks()`を使用してファイルを縮小する。

#### 【0046】

IFSによって、ファイル・ブロック関連の操作を行うために、下記の3つのインターフェースが使用される。

1. `csmAttachBlocks()`：IFSの意図が、`write()`である時に、このインターフェースが使用される。CSMは、そのキャッシュを介してこの要求を満足できない場合に、型`stpMsgType_BlkDiskAllocate`のトランザクションをサーバに送る。

2. `csmTranslateBlocks()`：このインターフェースは、`read()`または`write()`のいずれかに使用される。IFSは、読取中およびキャッシュのハードニング中（書込動作の一部）に、このインターフェースを使用することができる。CSMは、そのキャッシュを介してこの要求を満足できない場合に、型`stpMsgType_BlkDiskGetSegment`のトランザクションをサーバに送る。

3. `csmDetachBlocks()`：これは、`truncate()`に使用される。

#### 【0047】

本発明の態様によれば、2タイプの仮想 - 物理マッピングが、CSMで維持される。

1. 読取変換：これによって、読取用の仮想 - 物理マッピングが指定される。
2. 書込変換：これによって、書込用の仮想 - 物理マッピングが指定される。

#### 【0048】

一実施形態で、セグメントが、下記の3つの状態の読取および書込の変換リストを有する場合がある。

- ・有効な読取変換と無効な書込変換。
- ・有効な書込変換と無効な読取変換。
- ・読取変換と書込変換の両方が有効である。

10

20

30

40

50

## 【 0 0 4 9 】

## 読取

- ・読取システム呼出しについて、I F S は、読取フラグをセットして `csmTranslateBlocks()` を呼び出して、読取変換が必要であることを示す。
- ・「読取」変換について、C S M は、書込変換を先に調べる。書込変換が存在し、使用中である場合には、書込ブロック変換を返す。
- ・そうでない場合には、C S M は、読取変換を調べる。読取変換が使用可能であり、使用中である場合には、その読取変換が返される。
- ・そうでない場合には、C S M は、これらのブロックを 0 で充てんする必要があることを示す、0 を返す。

10

## 【 0 0 5 0 】

## 書込

- ・書込システム呼出しについて、I F S は、`csmAttachBlocks()` を呼び出す。成功の場合に、バッキング・ブロックが割り振られていることが保証される。
- ・更新の場合には、I F S は、そのブロックをキャッシュに戻す必要がある。したがって、PageInスレッドによって、読取フラグ付きで `csmTranslateBlocks()` を呼び出す。上の「読取」の論理から、I F S はブロック変換を得る（新規ブロックへの書込について、このステップをスキップすることができる）。
- ・ここで、I F S は、キャッシュ内ページを更新し、そのページをディスクにフラッシュする準備ができたならば、書込フラグ付きで `csmTranslateBlocks()` をもう一度呼び出す。ここで、C S M は、書込変換だけを与える必要がある。
- ・I F S は、前のステップで与えられた変換を使用し、キャッシュをフラッシュする（ディスクに書き込む）。

20

## 【 0 0 5 1 】

## 切捨

- ・I F S が、ファイル縮小のために `csmDetachBlocks()` を呼び出す。
- ・C S M は、対応する仮想ブロックの読取変換および書込変換の両方を、無効状態にマークする必要がある。

## 【 0 0 5 2 】

- 上の動作は、C S M のキャッシュで行われる。定期的な間隔で、または指定された状態で、C S M は、`blkdisk` 更新を介する修正によってサーバを更新することができる。また、上の説明では、主として、書込が「キャッシング付き / バッファ付き入出力」であると仮定される。「直接入出力」書込について、多少の差がある。

30

## 【 0 0 5 3 】

## 直接入出力

- 直接入出力は、必ずしもブロック境界で行われないので、`write()` の位置合せされない部分についてキャッシュ入出力を模倣する必要がある場合がある。

## 【 0 0 5 4 】

- `stpMsgType_BlkDiskAllocate` または `stpMsgType_BlkDiskGetSegment` のいずれかに応答して、C S M が、要求される各セグメントのエクステンツのリストを得る。これらのセグメント変換は、アンマーシャルされ、C S M のキャッシュに保管される。

40

## 【 0 0 5 5 】

- サーバは、書込変換と正確に同一である場合に読取変換を送らないようにすることができる。これが最適化になる可能性がある。

## 【 0 0 5 6 】

- 各セグメントは、`mcBlkDiskSegment` という名前のデータ構造体によって表される。

## 【 0 0 5 7 】

- このキャッシュ内セグメント構造体の要素に、下記を含めることができる。

## 【 0 0 5 8 】

【表 1】

s_objP	このセグメントが属するファイル・オブジェクトへのポインタ。
s_segNo	ファイル内のセグメント番号。
s_readExtentCount	このセグメント内の読取エクステントの現在の個数。
s_readExtentList	このセグメントの「読取」変換を表す連続するブロック・セグメントのリスト。
s_inlineReadExtent	上のリストのインライン表現。
s_readBlockUsedState	ビット・マップである。セグメント内のブロックごとに1ビット。ブロックが使用中／未使用状態のどちらであるかを示す。1－使用中、0－未使用。
s_writeExtentCount	このセグメントの書込エクステントの現在の個数。
s_writeExtentList	このセグメントの「書込」変換を表す連続するブロック・セグメントのリスト。
s_writeBlockUsedStat	ビット・マップである。セグメント内のブロックごとに1ビット。ブロックが使用中／未使用状態のどちらであるかを示す。1－使用中、0－未使用。
s_isDirty	ライブ・ブロック状態ビット・ベクトルに、サーバと同期化されなければならない更新が含まれる場合に真。
s_extentListValid	エクステント・リストが有効であり、キャッシュ内にある場合に真。

10

20

## 【0059】

クライアントは、排他データ・ロックの下でs\_readBlockUsedStateおよびs\_writeBlockUsedStateだけを変更する。セグメント変換の残りの部分は、クライアント側では変更されない。したがって、更新を送っている間に、クライアントは、この2つのビット・マップだけをサーバに送信する。

30

## 【0060】

読取動作では、ビット・マップが全く変更されない、すなわち、s\_readBlockUsedStateまたはs\_writeBlockUsedStateのいずれかが、読取変換をもたらすのに使用されるが、変更されない。

## 【0061】

書込動作では、s\_writeBlockUsedStateだけが操作され、s\_readBlockUsedStateは使用されない。読取動作とは異なって、s\_writeBlockUsedStateのビットマップを変更（セットのみ）して、成功裡の書込動作を示す場合がある。

40

## 【0062】

切捨動作では、両方のビット・マップが変更される可能性がある。

## 【0063】

したがって、短く言えば、読取動作では何も変更されず、書込動作ではs\_writeBlockUsedStateビットマップ・ベクトルの少数のビットがセットされる場合があり、切捨動作では両方のビットマップ・ベクトルのビットがアンセットされる場合がある。

## 【0064】

要約すると、キャッシュ付き入出力の場合に、コピー・オン・ライト（COW）を、ほとんどコストなしで達成することができる。これは、通常、書込動作の完了にかかわる下記の2つの異なるスレッド／動作があるからである。

50

## 【 0 0 6 5 】

1 . ターゲット・データをキャッシュ内で更新 / 変更させるPageInスレッド / 動作と、

## 【 0 0 6 6 】

2 . 更新されたページをディスクにフラッシュ・バックするPageOutスレッド / 動作。

## 【 0 0 6 7 】

これを与えられて、C O Wを、PageInに読取変換、PageOutに書込変換を使用することによって実行することができる。

## 【 0 0 6 8 】

例として下記を検討されたい。

## 【 0 0 6 9 】

新しいファイルの作成の際に、C S Mは、

- 書込ブロック・エクステンツ
- N U L L 読取ブロック・エクステンツ

を得る。

## 【 0 0 7 0 】

新しいファイルなので、ページ・インすべきものではなく、したがって、書込データは、キャッシュの空白ページに送られる。

## 【 0 0 7 1 】

書込が進行する際に、I F Sは、( C S Mインターフェースを介して ) s\_writeBlockUsedStateビット・ベクトルの対応するビットをセットし、それらが使用中であることを示す。

## 【 0 0 7 2 】

上で説明したように、このブロック範囲内の将来の読取および書込では、書込エクステンツからの変換を得る。

## 【 0 0 7 3 】

管理者が、フラッシュ・コピーをとると仮定する。

## 【 0 0 7 4 】

フラッシュ・コピー動作の一部として、サーバによって、クライアントからのすべてのデータ・ロックが取り消される。したがって、修正されたデータが、ディスクに同期化され、メタデータが、更新トランザクションを介してサーバに送られる。クライアントは、データ・ロックを有しないので、その変換のすべてが無効になる。

## 【 0 0 7 5 】

フラッシュ・コピーの後

読取について、クライアントは、csmTranslateBlocks()を呼び出し、サーバは、「読取」変換を呼び出すことができるが、「書込」変換はN U L Lになる。

## 【 0 0 7 6 】

書込について、クライアントは、csmAttachBlocks()を呼び出す。サーバは、読取変換を返し(上と同一)、書込エクステンツ・リストについて、サーバは、未使用ブロックの新しい組を割り振り、返さなければならない。したがって、クライアントは、2つの変換を有する。クライアントは、「ページイン」の一部として読取変換を使用し、「ページアウト」の一部として書込変換を使用する。

## 【 0 0 7 7 】

したがって、フラッシュ・コピーの一部であったデータ・ブロックを、バッファ・キャッシュにPageInし、更新を適用する。ページ・アウト中に、書込変換を使用するが、この書込は、新しい物理ブロックへのPageOutスレッドをポイントする。

## 【 0 0 7 8 】

directI0のC O Wは、多少異なる。直接入出力について：

o I F Sは、書込フラグをセットしてcsmTranslateBlocks()を呼び出して、書込を進めるためのバッキング・ブロックを有するかどうかを調べる。C S Mからの変換を分析した後、バッキング・ブロックが割り振られていない場合に、I F Sは、csmAttachBlock

10

20

30

40

50

s()を呼び出す。

- o csmAttachBlocks()について、I F Sは、そのデータロックをS H A R E D \_ W R I T EモードからE X C L U S I V Eに切り替える必要がある。
- o 書込が終了した後に、クライアントは、これらのブロックをU S E D状態にマークする。この時に、クライアントは、排他モードのデータ・ロックを保持する必要がある。
- o 直接入出力境界が、ブロック・サイズに位置合わせされていない場合に、最初のブロックおよび最後のブロックへの入出力が、キャッシング付きの形で行われる。これを行うステップに、下記が含まれる。
  - ・「読取変換」についてC S Mを呼び出す。
  - ・ローカル・カーネル・バッファを割り振る。
  - ・新たに割り振られたカーネル・バッファにディスク・ブロックを読み込む。
  - ・カーネル・バッファを更新する。
  - ・書込変換についてC S Mを呼び出す。
  - ・新しいブロック位置に書き込む（書込変換）
  - ・C S Mを呼び出して新しい書込ブロックの「U S E D」ビットをマークする。

10

#### 【0079】

中間ブロックについて、

- 境界をブロック・サイズに位置合せする。
- C S Mを呼び出して、書込変換を得る。
- ディスクに書き込む。
- C S Mを呼び出して、ブロックの「U S E D」ビットをマークする。

20

#### 【0080】

下記は、クライアントでのビット・メトリックスの例である。

#### 【0081】

WriteBlockUsedStateビット配列：Wbit

#### 【0082】

ReadBlockUsedStateビット配列：Rbit

#### 【0083】

かぎ括弧内で、( @ < > )は、物理ブロック・アドレスを示す。すなわち( @ A )は、物理ブロック・アドレスがAであることを意味する。

30

#### 【0084】

U D - 未定義

#### 【0085】

【表 2】

	仮想ブロック#1	仮想ブロック#2	仮想ブロック#3	仮想ブロック#4				
オペランド	(W Phy Blk @) Wbit	(R Phy Blk @) Rbit	(W Phy Blk @) Wbit	(R Phy Blk @) Rbit	(W Phy Blk @) Wbit	(R Phy Blk @) Rbit	(W Phy Blk @) Wbit	(R Phy Blk @) Rbit
ファイル作 成および割 振	(@A) 0	UD UD	(@B) 0	UD UD	(@C) 0	UD UD	(@D) 0	UD UD
3 ブロック 書込	(@A) 1	UD UD	(@B) 1	UD UD	(@C) 1	UD UD	(@D) 0	UD UD
ファイル読 取(上から変 更なし)	(@A) 1	UD UD	(@B) 1	UD UD	(@C) 1	UD UD	(@D) 0	UD UD
2 ブロック に切捨	(@A) 1	UD UD	(@B) 1	UD UD	(@C) 0	UD UD	(@D) 0	UD UD
読取	NULL UD	(@ A) 1	NULL UD	(@B) 1	UD UD	NULL UD	UD UD	NULL UD
最初の 2 ブ ロックの更 新を希望	(@X) 0	(@A) 1	(@Y) 0	(@B) 1	UD UD	NULL UD	UD UD	NULL UD
最初の 2 ブ ロックを書 込 (更新)	(@X) 1	(@A) 1	(@Y) 1	(@B) 1	UD UD	NULL UD	UD UD	NULL UD
1 ブロック に切捨	(@X) 1	(@A) 1	(@Y) 0	(@B) 0	UD UD	NULL UD	UD UD	NULL UD

10

20

30

## 【0086】

クライアントは、クライアント・サーバ環境で、WriteBlockUsedStateビット配列およびReadBlockUsedStateビット配列だけをサーバに送信する。

## 【0087】

C S Mは、これらのビットを解釈し、下記のようにマッピングする

40

## 【0088】



【表 3】

X = このブロックに関するマッピングなし

M = このブロックに関するマッピングが存在する

状態	書込マッピング	書込ビット	読取マッピング	読取ビット
未割振り	X	0	X	0
未定義	X	0	X	1
未定義	X	0	M	0
読取可能、共用 (COW ブロック割振りが必要)	X	0	M	1
未定義	X	1	X	0
未定義	X	1	X	1
未定義	X	1	M	0
未定義	X	1	M	1
書込可能だが読取不能	M	0	X	0
未定義	M	0	X	1
未定義	M	0	M	0
COW 保留中、読取で RM、書込で WM を使用	M	0	M	1
書込マッピングを介して書込可能かつ読取可能	M	1	X	0
未定義	M	1	X	1
未定義	M	1	M	0
未定義	M	1	M	1

10

20

## 【0089】

ブロックのマッピングを判定するために、CSMは、まず、書込マッピングを調べ、書込マッピングが存在し、対応するWビットがセットされている場合に、そのマッピングを、読取動作および書込動作の両方に使用する。

30

## 【0090】

書込マッピングが存在するが、Wビットが0である場合には、CSMは、読取マッピングを探す。

## 【0091】

読取マッピングが存在しない場合には、そのブロックは未使用（未初期化）とみなされ、そのブロックは、書込動作だけに使用することができる。すべての読取動作によって、バッファに0が充てんされる。

40

## 【0092】

読取マッピングが存在する場合には、そのブロックは、使用中（初期化済み）とみなされ、CSMは、Rビットが1であると仮定する（そうでなければならない）。この状態のブロックは、COW保留中とのみみなされ、その内容を修正するにはCOWが必要である。

## 【0093】

書込マッピングが存在しない場合には、CSMは、読取マッピングを探す。読取マッピングが存在する場合に、CSMは、Rビットが1であると仮定し（そうでなければならない）、そのブロックは、読取動作に使用中（初期化済み）だが書込動作には使用されていないとみなされる。書込動作を実行する前に、CMSは、サーバによって新しいバッキン

50

グ・ブロックが割り振られることを要求しなければならない。

【 0 0 9 4 】

C S Mは、サーバにビット・ベクトルを返して、使用中、切捨済み、またはコピー・オン・ライト済みあるいはこれらの組合せとしてブロックの状況の変化を示す。リターンの際に、サーバは、下記のようにビットを解釈する。

【 0 0 9 5 】

【表 4】

X = ドントケア

状態	W	R
ブロックはライブである（COWまたは使用中の割り振られたブロックをハードニング）	1	X
このブロックを切り捨てる（使用中の場合に未使用としてマークされる－解放可能）	0	0
このブロックに対する変更なし（COW保留を取り消すことができる）	0	1

10

【 0 0 9 6 】

書込ビットがセットされている場合には：

20

サーバは、読取ビットを無視する。

ブロックがAllocated状態であった（すなわち、従来の形で割り振られたがまだライブでない）場合には、Live状態に変更される（すなわち、ブロックが書込可能かつ読取可能）。

ブロックがCOW\_Pending状態であった（すなわち、ブロックが、読取マッピングおよび書込マッピングの両方を有し、これらがCOWにおいて異なる）場合には、Live状態になる（すなわち、ブロックが、書込マッピングであったものを介して書込可能かつ読取可能）。

【 0 0 9 7 】

書込ビットがセットされず、読取ビットがセットされる場合には：

30

これによって、問題のブロックに関する変更がないことがサーバに示される。

ブロックが、Allocated、COW\_Pending、またはPIT\_COW\_Pending状態であった場合には、そのブロックは、この状態に留まることができ、あるいは、そのブロックを非同期に解放することができる。クライアントは、そのブロックが解放されることまたはCOW\_Pending状態の1つに留まるかどうかを仮定することができない。

【 0 0 9 8 】

書込ビットがセットされず、読取ビットがセットされない場合には：

これによって、ブロックが切り捨てられた（前に割り振られたと仮定して）ことがサーバに示される。

ブロックがUnallocated状態であった場合には、そのブロックはUnallocated状態のままになる。 40

ブロックがAllocated状態であった場合には、そのブロックはAllocated状態のままになる。

ブロックがLive状態であった場合には、そのブロックはAllocated状態になる。

ブロックがShared状態であった場合には、そのブロックはUnallocated状態になり、読取（専用）マッピングが破棄される。

ブロックがCOW\_Pending状態であった場合には、そのブロックはAllocated状態になる。

ブロックがPIT\_COW\_Pending状態であった場合には、そのブロックはAllocated状態になる。

【 0 0 9 9 】

50

長所

有利なことに、本明細書で、冗長な入出力を最小限にすることによって、従来のデータ・ファイル書込処理に対する最小限の追加コストで、コンピューティング環境内でコピー・オン・ライトを実施する技法が提示される。提示されたコピー・オン・ライト技法は、標準ファイルシステム・ドライバなどの上位層に透過的である。この技法には、2つの異なる同時変換すなわち、読取マッピング・テーブルおよび書込マッピング・テーブルを使用して、単一の書込動作を使用してファイル内のデータの単位のコピー・オン・ライトを達成することが含まれる。もう1つの態様では、クライアント・サーバ環境の複数のクライアントにまたがるファイルの分散コピー・オン・ライトを実施する技法が提示される。有利なことに、この分散コピー・オン・ライト実施形態では、中央サーバへの負荷が減り、クライアントの追加に伴ってスケールアップされ、1つまたは複数のクライアントが動作不能になる場合があるにもかかわらず、ファイルのコピー・オン・ライトを進行できるようになる。さらに、提示された分散コピー・オン・ライトでは、並列コピー・オン・ライトが可能であり、作業負荷を複数のクライアントの間で分散でき、これによって、リソースのより効率的な利用がもたらされる。

【0100】

代替実施形態

コンピューティング環境の例を提供したが、これらは例にすぎない。他の実施形態を使用することができる。たとえば、本明細書で、ファイルシステムに関する例を説明したが、これは単なる一例にすぎない。本発明の1つまたは複数の他の態様は、他の環境に適用可能である。

【0101】

本発明を、たとえばコンピュータ使用可能媒体を有する製造品（たとえば、1つまたは複数のコンピュータ・プログラム製品）に含めることができる。この媒体は、その中に、たとえば、本発明の機能を提供し促進するコンピュータ可読プログラム・コード手段または論理（たとえば、命令、コード、コマンドなど）を実施される。製造品を、コンピュータ・システムの一部として含めるか、別々に販売することができる。

【0102】

さらに、本発明の機能を実行するために機械によって実行可能な命令の少なくとも1つのプログラムを実施する、機械によって可読の少なくとも1つのプログラム記憶装置を提供することができる。

【0103】

本明細書で示された流れ図は、例にすぎない。本発明の趣旨から逸脱しない、これらの図または図に示されたステップ（または動作）に対する多数の変形形態がありえる。たとえば、ステップを、異なる順序で実行することができ、ステップを追加、削除、または修正することができる。これらの変形形態のすべてが、請求される発明の一部とみなされる。

【0104】

好ましい実施形態を図示し、本明細書で詳細に説明したが、本発明の趣旨から逸脱せずに、さまざまな修正、追加、置換、および類似物を行うことができ、したがって、これらが、請求項で定義される発明の趣旨に含まれるとみなされることを、当業者は諒解するであろう。

【0105】

まとめとして、本発明の構成に関して以下の事項を開示する。

【0106】

- (1) コンピューティング環境でコピー・オン・ライトを実施する方法であって、
  - (i) ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第1仮想ブロック・物理ブロック・マッピングを実行するのに第1マッピング・テーブルを使用することと、
  - (ii) 前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むの

に使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに第2マッピング・テーブルを使用することであって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、使用することを含む方法。

(2) 前記第1マッピング・テーブルが、読取マッピング・テーブルを含み、前記第2マッピング・テーブルが、書込マッピング・テーブルを含み、前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック・物理ブロック・マッピングと異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも1つの仮想ブロックを含む、上記(1)に記載の方法。

(3) 前記コピー・オン・ライト実施形態が、さらに、まず修正が部分ブロック書込または全ブロック書込のどちらを含むかを判定することと、部分ブロック書込の場合に、前記使用すること(i)および前記使用すること(ii)を実行し、そうでない場合に前記使用すること(i)を実行せずに前記実行すること(ii)を実行することを含む、上記(1)に記載の方法。 10

(4) 前記実行すること(i)が、データの前記ブロックを物理記憶装置からバッファに読み取ることを含み、前記方法が、さらに、前記使用すること(ii)を実行する前に前記バッファ内のデータの前記ブロックを修正することを含む、上記(1)に記載の方法。

(5) 前記コンピューティング環境が、ファイルシステム・サーバおよび少なくとも1つのクライアントを含むクライアント・サーバ環境を含み、前記使用すること(i)および前記使用すること(ii)が、前記クライアント・サーバ環境の前記少なくとも1つのクライアントによって実行される、上記(1)に記載の方法。 20

(6) 前記少なくとも1つのクライアントによって実行される前記使用すること(i)および前記使用すること(ii)が、前記ファイルのデータの前記ブロックのコピー・オン・ライトを実行する時に、前記第1マッピング・テーブルおよび前記第2マッピング・テーブルの少なくとも1つを入手するために前記ファイルシステム・サーバへの少なくとも1つの呼出しを行うことをさらに含む、上記(5)に記載の方法。

(7) データの前記修正されたブロックを物理記憶装置に書き込んだ後に、前記第1マッピング・テーブルを更新することをさらに含み、前記更新することが、前記第1マッピング・テーブルの少なくとも1つの仮想ブロック・物理ブロック変換を、前記第2マッピング・テーブルの対応する仮想ブロック・物理ブロック変換と一致するように修正することを含む、上記(5)に記載の方法。 30

(8) 前記コンピューティング環境が、コンピューティング・ユニットおよび外部記憶ユニットを含み、前記外部記憶ユニットが、前記物理記憶装置を含み、前記使用すること(i)および前記使用すること(ii)が、前記コンピューティング・ユニットによって実行される、上記(1)に記載の方法。

(9) クライアント・サーバ・コンピューティング環境でコピー・オン・ライトを容易にする方法であって、

前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持すること 40

を含み、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第1仮想ブロック・物理ブロック・マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能である方法。

(10) 前記ファイルに関する前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック・物理ブロック変換と異なる、前記物理記憶装置の物理 50

ブロックにマッピングされる少なくとも1つの仮想ブロックを含む、上記(9)に記載の方法。

(11) 前記ファイルのデータのブロックのコピー・オン・ライトが実行された後に前記読取マッピング・テーブルを更新することをさらに含み、前記更新することが、前記読取マッピング・テーブルの少なくとも1つの仮想ブロック - 物理ブロック変換を、前記書込マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と一致するように修正することを含む、上記(9)に記載の方法。

(12) 複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施する方法であって、

前記クライアント・サーバ環境の複数のクライアントを使用してファイルのコピー・オン・ライトを実行することであって、 10

(i) 前記複数のクライアントの第1クライアントによって、前記ファイルのデータの少なくとも1つのブロックのコピー・オン・ライトを実行することと、

(ii) 前記複数のクライアントの第2クライアントによって、前記ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトを実行することと

を含む、実行すること

を含む方法。

(13) 前記実行すること(i)が、前記第1クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行することを含み、前記実行すること(ii)が、前記第2クライアントによって、 20  
単一の書込動作を使用してデータの前記少なくとも1つの他のブロックのコピー・オン・ライトを実行することを含む、上記(12)に記載の方法。

(14) 前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記ファイルシステム・サーバが、前記少なくとも1つの共用記憶装置に保管された前記ファイルに関するマッピング・テーブルを維持し、前記実行すること(i)が、前記第1クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを入手することと、前記ファイルのデータの前記少なくとも1つのブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用することとを含み、前記実行すること(ii)が、前記第 30  
2クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを入手することと、前記ファイルのデータの前記少なくとも1つの他のブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用することとを含む、上記(13)に記載の方法。

(15) 前記第1クライアントによって、前記第1クライアントが前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行したことを前記ファイルシステム・サーバに知らせることと、それに応答して、前記ファイルシステム・サーバによって維持される前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも1つを更新することとをさらに含み、上記(14)に記 40  
載の方法。

(16) 前記ファイルシステム・サーバが、前記第1クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも1つのブロックおよび前記第2クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも1つの他のブロックに対するコピー・オン・ライトを前記複数のクライアントのすべてのクライアントが行えなくする、上記(14)に記載の方法。

(17) 前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記方法が、さらに、前記ファイルの前記コピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファ 50  
イルを制御することを含み、前記開始することが、前記実行すること(i)および前記実

行すること ( i i ) によって使用される前記ファイルに関する書込マッピング・テーブルを更新することを含む、上記 ( 1 2 ) に記載の方法。

( 1 8 ) クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にする方法であって、

前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御することであって、前記制御することが、前記クライアント・サーバ環境の第 1 クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにすることと、前記クライアント・サーバ環境の第 2 クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにすることとを含み、前記ファイルシステム・サーバが、前記ファイルの分散コピー・オン・ライトの実行を制御し、容易にする、制御すること

10

を含む方法。

( 1 9 ) 前記制御することが、前記ファイルのコピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御することを含み、前記開始することが、前記コピー・オン・ライトを実行するのに使用される前記ファイルに関する少なくとも 1 つのマッピング・テーブルを更新することを含む、上記 ( 1 8 ) に記載の方法。

( 2 0 ) 前記ファイルシステム・サーバで、前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持することをさらに含み、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルが、前記コピー・オン・ライトの実行に使用される、上記 ( 1 8 ) に記載の方法。

20

( 2 1 ) 前記コピー・オン・ライトの実行の後に、前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも 1 つを更新することをさらに含む、上記 ( 2 0 ) に記載の方法。

( 2 2 ) 前記制御することが、前記ファイルシステム・サーバによって、前記ファイルの前記コピー・オン・ライトの一部としての、前記第 1 クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記部分に対する追加更新および前記第 2 クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記異なる部分に対する追加更新を防ぐことをさらに含む、上記 ( 1 8 ) に記載の方法。

( 2 3 ) コンピューティング環境でコピー・オン・ライトを実施するシステムであって、

30

( i ) ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第 1 仮想ブロック - 物理ブロック・マッピングを実行するのに第 1 マッピング・テーブルを使用する手段と、

( i i ) 前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第 2 仮想ブロック - 物理ブロック・マッピングを実行するのに第 2 マッピング・テーブルを使用する手段であって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、使用する手段と

を含むシステム。

( 2 4 ) 前記第 1 マッピング・テーブルが、読取マッピング・テーブルを含み、前記第 2 マッピング・テーブルが、書込マッピング・テーブルを含み、前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック - 物理ブロック・マッピングと異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも 1 つの仮想ブロックを含む、上記 ( 2 3 ) に記載のシステム。

40

( 2 5 ) 前記コピー・オン・ライト実施形態が、さらに、まず修正が部分ブロック書込または全ブロック書込のどちらを含むかを判定し、部分ブロック書込の場合に、前記使用すること ( i ) および前記使用すること ( i i ) を実行し、そうでない場合に前記使用すること ( i ) を実行せずに前記実行すること ( i i ) を実行する手段を含む、上記 ( 2 3 ) に記載のシステム。

( 2 6 ) 前記実行する手段 ( i ) が、データの前記ブロックを物理記憶装置からバッファに読み取る手段を含み、前記システムが、さらに、前記使用すること ( i i ) を実行する

50

前に前記バッファ内の前記ブロックを修正する手段を含む、上記(23)に記載のシステム。

(27) 前記コンピューティング環境が、ファイルシステム・サーバおよび少なくとも1つのクライアントを含むクライアント・サーバ環境を含み、前記使用する手段(i)および前記使用する手段(ii)が、前記クライアント・サーバ環境の前記少なくとも1つのクライアントによって実行される、上記(23)に記載のシステム。

(28) 前記少なくとも1つのクライアントによって実行される前記使用する手段(i)および前記使用する手段(ii)が、前記ファイルのデータの前記ブロックのコピー・オン・ライトを実行する時に、前記第1マッピング・テーブルおよび前記第2マッピング・テーブルの少なくとも1つを入手するために前記ファイルシステム・サーバへの少なくとも1つの呼出しを行う手段をさらに含む、上記(27)に記載のシステム。 10

(29) データの前記修正されたブロックを物理記憶装置に書き込んだ後に、前記第1マッピング・テーブルを更新する手段をさらに含み、前記更新する手段が、前記第1マッピング・テーブルの少なくとも1つの仮想ブロック - 物理ブロック変換を、前記第2マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と一致するように修正する手段を含む、上記(27)に記載のシステム。

(30) 前記コンピューティング環境が、コンピューティング・ユニットおよび外部記憶ユニットを含み、前記外部記憶ユニットが、前記物理記憶装置を含み、前記使用する手段(i)および前記使用する手段(ii)が、前記コンピューティング・ユニットによって実行される、上記(23)に記載のシステム。 20

(31) クライアント・サーバ・コンピュータ環境でコピー・オン・ライトを容易にするシステムであって、

前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する手段

を含み、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第1仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック - 物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能である 30

システム。

(32) 前記ファイルに関する前記書込マッピング・テーブルが、前記読取マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と異なる、前記物理記憶装置の物理ブロックにマッピングされる少なくとも1つの仮想ブロックを含む、上記(31)に記載のシステム。

(33) 前記ファイルのデータのブロックのコピー・オン・ライトが実行された後に前記読取マッピング・テーブルを更新する手段をさらに含み、前記更新する手段が、前記読取マッピング・テーブルの少なくとも1つの仮想ブロック - 物理ブロック変換を、前記書込マッピング・テーブルの対応する仮想ブロック - 物理ブロック変換と一致するように修正する手段を含む、上記(31)に記載のシステム。 40

(34) 複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施するシステムであって、

(i) 前記クライアント・サーバ環境の第1クライアントで、コピー・オン・ライトされる前記ファイルのデータの少なくとも1つのブロックのコピー・オン・ライトを実行する手段と、

(ii) 前記クライアント・サーバ環境の第2クライアントで、コピー・オン・ライトされる前記ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトを実行する手段であって、前記ファイルの前記コピー・オン・ライトの異なる部分が、前記 50

クライアント・サーバ環境内の前記複数のクライアントの異なるクライアントによって実行される、手段と

を含むシステム。

(35) 前記実行する手段(i)が、前記第1クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行する手段を含み、前記実行する手段(ii)が、前記第2クライアントによって、単一の書込動作を使用してデータの前記少なくとも1つの他のブロックのコピー・オン・ライトを実行する手段を含む、上記(34)に記載のシステム。

(36) 前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記ファイルシステム・サーバが、前記少なくとも1つの共用記憶装置に保管された前記ファイルに関するマッピング・テーブルを維持し、前記実行する手段(i)が、前記第1クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも1つのブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する手段を含み、前記実行する手段(ii)が、前記第2クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも1つの他のブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する手段を含む、上記(35)に記載のシステム。 10 20

(37) 前記第1クライアントによって、前記第1クライアントが前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行したことを前記ファイルシステム・サーバに知らせ、それに応答して、前記ファイルシステム・サーバによって維持される前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも1つを更新する手段をさらに含む、上記(36)に記載のシステム。

(38) 前記ファイルシステム・サーバが、前記第1クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも1つのブロックおよび前記第2クライアントによってコピー・オン・ライト更新された前記ファイルのデータの前記少なくとも1つの他のブロックに対するコピー・オン・ライトを前記複数のクライアントのいずれもが行えなくする、上記(36)に記載のシステム。 30

(39) 前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記システムが、さらに、前記ファイルの前記コピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御する手段を含み、前記開始する手段が、前記実行する手段(i)および前記実行する手段(ii)によって使用される前記ファイルに関する書込マッピング・テーブルを更新する手段を含む、上記(34)に記載のシステム。

(40) クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にするシステムであって、 40

前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御する手段であって、前記制御する手段が、前記クライアント・サーバ環境の第1クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにし、前記クライアント・サーバ環境の第2クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにする手段を含み、前記ファイルシステム・サーバが、前記ファイルの分散コピー・オン・ライトの実行を制御し、容易にする、制御する手段

を含むシステム。

(41) 前記制御する手段が、前記ファイルのコピー・オン・ライトを開始するために前記ファイルシステム・サーバによって前記ファイルを制御する手段を含み、前記開始する 50



ことが、前記コピー・オン・ライトを実行するのに使用される前記ファイルに関する少なくとも1つのマッピング・テーブルを更新することを含む、上記(40)に記載のシステム。

(42) 前記ファイルシステム・サーバで、前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する手段をさらに含み、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルが、前記コピー・オン・ライトの実行に使用される、上記(40)に記載のシステム。

(43) 前記コピー・オン・ライトの実行の後に、前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルの少なくとも1つを更新する手段をさらに含み、上記(42)に記載のシステム。

10

(44) 前記制御する手段が、前記ファイルシステム・サーバによって、前記ファイルの前記コピー・オン・ライトの一部としての、前記第1クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記部分に対する追加更新および前記第2クライアントによってコピー・オン・ライト更新された前記ファイル内の前記データの前記異なる部分に対する追加更新を防ぐ手段をさらに含み、上記(40)に記載のシステム。

(45) 製造品であって、

コンピューティング環境でコピー・オン・ライトを実施するコンピュータ可読プログラム・コード論理を有する少なくとも1つのコンピュータ使用可能媒体

20

を含み、前記コンピュータ可読プログラム・コード論理が、

(i) ファイルのデータのブロックを修正のために物理記憶装置から読み取るのに使用される第1仮想ブロック - 物理ブロック・マッピングを実行するのに第1マッピング・テーブルを使用する論理と、

(ii) 前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック - 物理ブロック・マッピングを実行するのに第2マッピング・テーブルを使用する論理であって、データの前記ブロックのコピー・オン・ライトが、単一の書込動作を使用して達成される、使用する論理と

を含む、製造品。

(46) 前記コンピューティング環境が、ファイルシステム・サーバおよび少なくとも1つのクライアントを含むクライアント・サーバ環境を含み、前記使用する論理(i)および前記使用する論理(ii)が、前記クライアント・サーバ環境の前記少なくとも1つのクライアントによって実行される、上記(45)に記載の製造品。

30

(47) 前記少なくとも1つのクライアントによって実行される前記使用する論理(i)および前記使用する論理(ii)が、前記ファイルのデータの前記ブロックのコピー・オン・ライトを実行する時に、前記第1マッピング・テーブルおよび前記第2マッピング・テーブルの少なくとも1つを入手するために前記ファイルシステム・サーバへの少なくとも1つの呼出しを行うことをさらに含み、上記(45)に記載の製造品。

(48) 前記コンピューティング環境が、コンピューティング・ユニットおよび外部記憶ユニットを含み、前記外部記憶ユニットが、前記物理記憶装置を含み、前記使用する論理(i)および前記使用する論理(ii)が、前記コンピューティング・ユニットによって

40

(49) 製造品であって、

クライアント・サーバ・コンピューティング環境でコピー・オン・ライトを容易にするコンピュータ可読プログラム・コード論理を有する少なくとも1つのコンピュータ使用可能媒体

を含み、前記コンピュータ可読プログラム・コード論理が、

前記クライアント・サーバ・コンピューティング環境のファイルシステム・サーバで、ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する論理であって、前記読取マッピング・テーブルが、修正のために物理記憶装置から前記ファイルのデータのブロックを読み取るのに使用される第1仮想ブロック - 物理ブロック・

50

マッピングを実行するのに使用可能であり、前記書込マッピング・テーブルが、前記ファイルの前記データの修正されたブロックを物理記憶装置に書き込むのに使用される第2仮想ブロック・物理ブロック・マッピングを実行するのに使用可能であり、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用して、データのブロックのコピー・オン・ライトが単一の書込動作を使用して達成可能である、論理

を含む、製造品。

(50) 製造品であって、

複数のクライアントを有するクライアント・サーバ環境内でファイルのコピー・オン・ライトを実施するコンピュータ可読プログラム・コード手段を有する少なくとも1つのコンピュータ使用可能媒体

10

を含み、前記コンピュータ可読プログラム・コード手段が、

前記クライアント・サーバ環境の複数のクライアントを使用してファイルのコピー・オン・ライトを実行する論理

を含み、前記実行する論理が、

(i) 前記複数のクライアントの第1クライアントによって、前記ファイルのデータの少なくとも1つのブロックのコピー・オン・ライトを実行する論理と、

(ii) 前記複数のクライアントの第2クライアントによって、前記ファイルのデータの少なくとも1つの他のブロックのコピー・オン・ライトを実行する論理と

を含む、製造品。

(51) 前記論理(i)が、前記第1クライアントによって、単一の書込動作を使用して前記ファイルのデータの前記少なくとも1つのブロックのコピー・オン・ライトを実行する論理を含み、前記論理(ii)が、前記第2クライアントによって、単一の書込動作を使用してデータの前記少なくとも1つの他のブロックのコピー・オン・ライトを実行する論理を含む、上記(50)に記載の製造品。

20

(52) 前記クライアント・サーバ環境のファイルシステム・サーバが、前記ファイルを含む少なくとも1つの共用記憶装置に関連し、前記ファイルシステム・サーバが、前記少なくとも1つの共用記憶装置に保管された前記ファイルに関するマッピング・テーブルを維持し、前記論理(i)が、前記第1クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも1つのブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する論理を含み、前記論理(ii)が、前記第2クライアントによって、前記ファイルシステム・サーバから前記ファイルに関する前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを入手し、前記ファイルのデータの前記少なくとも1つの他のブロックの前記コピー・アンド・ライトを実行する際に前記読取マッピング・テーブルおよび前記書込マッピング・テーブルを使用する論理を含む、上記(51)に記載の製造品。

30

(53) 製造品であって、

クライアント・サーバ環境内でファイルのコピー・オン・ライトを容易にするコンピュータ可読プログラム・コード論理を有する少なくとも1つのコンピュータ使用可能媒体

40

を含み、前記コンピュータ可読プログラム・コード論理が、

前記クライアント・サーバ環境の共用記憶ユニットに保管されたファイルのコピー・オン・ライトの実施をファイルシステム・サーバから制御する論理であって、前記制御する論理が、前記クライアント・サーバ環境の第1クライアントが前記ファイル内のデータの部分をコピー・オン・ライトできるようにする論理と、前記クライアント・サーバ環境の第2クライアントが前記ファイル内の前記データの異なる部分をコピー・オン・ライトできるようにする論理とを含み、前記ファイルシステム・サーバが、前記ファイルの分散コピー・オン・ライトの実行を制御し、容易にする、制御する論理

を含む、製造品。

(54) 前記制御する論理が、前記ファイルのコピー・オン・ライトを開始するために前

50

記ファイルシステム・サーバによって前記ファイルを制御する論理を含み、前記開始することが、前記コピー・オン・ライトを実行するのに使用される前記ファイルに関する少なくとも1つのマッピング・テーブルを更新することを含む、上記(53)に記載の製造品。

(55) 前記ファイルシステム・サーバで、前記ファイルに関する読取マッピング・テーブルおよび書込マッピング・テーブルを維持する論理をさらに含み、前記読取マッピング・テーブルおよび前記書込マッピング・テーブルが、前記コピー・オン・ライトの実行に使用される、上記(53)に記載の方法。

【図面の簡単な説明】

【0107】

10

【図1】本発明の1つまたは複数の態様が組み込まれ、使用されるコンピューティング環境の一実施形態を示す図である。

【図2】範囲またはオフセットが仮想/相対ブロック番号に変換される、アプリケーションまたはプロセスに対するファイル表現(200)の一例を示す図である。

【図3】図2の仮想/相対ブロック番号が、ファイルシステムの1つまたは複数の記憶ユニット内の物理ブロック・アドレスにマッピングされる、ファイルシステム・マッピング・テーブル(300)を示す図である。

【図4】ファイルシステムの記憶ユニット内の図3の物理ブロック・アドレスの配置を示す図である。

【図5】オフセット5000から開始され、ファイルシステムの記憶ユニットの物理ブロックD内に含まれる4バイトのデータの読取の例を示す図である。 20

【図6】本発明の態様による書込動作およびコピー・オン・ライト動作の一実施形態を示す流れ図である。

【図7】図6の論理に従う書込動作で使用されるデータのブロックを記憶ユニットから論理バッファまたはキャッシュに読み取る読取処理の例を示す流れ図である。

【図8】図6の論理に従う書込動作で使用されるデータの修正されたブロックを記憶ユニットに書き込む書込処理の例を示す流れ図である。

【図9】本発明の態様による、ファイルのコピー・オン・ライトに使用される読取マッピング・テーブル(900)および書込マッピング・テーブル(910)の例を示す図である。 30

【図10】本発明の態様による、図6の論理によるコピー・オン・ライト動作に使用されるデータ・ブロック読取処理の一実施形態を示す流れ図である。

【図11】本発明の態様による、図6の論理によるコピー・オン・ライト動作に関する修正済みデータ・ブロック書込処理の一実施形態を示す流れ図である。

【図12】本発明の態様による、ファイルのコピー・オン・ライトに使用される読取マッピング・テーブル(1200)および書込マッピング・テーブル(1210)のもう1つの例を示す図である。

【図13】本発明の1つまたは複数の態様が組み込まれ、使用されるコンピューティング環境のもう1つの例を示す図である。

【図14】本発明の態様による、たとえば図13のクライアント・サーバ環境の、ファイルの少なくとも1ブロックのデータのコピー・オン・ライトを実行するクライアント1の一実施形態を示す流れ図である。 40

【図15】本発明の態様による、たとえば図13のクライアント・サーバ環境の、ファイルの少なくとも1つの他のブロックのデータのコピー・オン・ライトを実行するクライアント2の一実施形態を示す流れ図である。

【図16】本発明の態様による、コピー・オン・ライト中に使用されるファイルに関する読取マッピング・テーブル(1600)および書込マッピング・テーブル(1610)のもう1つの例を示す図である。

【図17】本発明の態様による、コピー・オン・ライトを使用するデータ保存アプリケーションの一例を示す流れ図である。 50

【符号の説明】

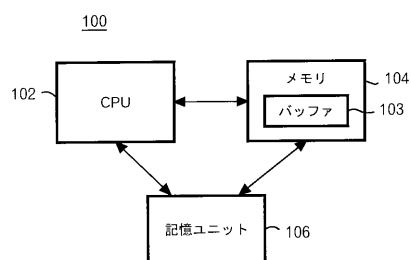
【 0 1 0 8 】

- |         |                     |
|---------|---------------------|
| 1 0 0   | コンピューティング環境         |
| 1 0 2   | 中央処理装置              |
| 1 0 4   | メモリ                 |
| 1 0 3   | バッファまたはキャッシュ領域      |
| 1 0 6   | 記憶ユニットまたは記憶装置       |
| 2 0 0   | ファイル表現              |
| 3 0 0   | ファイルシステム・マッピング・テーブル |
| 4 0 0   | 記憶ユニット              |
| 9 0 0   | 読取マッピング・テーブル        |
| 9 1 0   | 書込マッピング・テーブル        |
| 1 2 0 0 | 更新された読取マッピング・テーブル   |
| 1 2 1 0 | 書込マッピング・テーブル        |
| 1 3 0 0 | 環境                  |
| 1 3 1 0 | クライアント 1            |
| 1 3 2 0 | クライアント 2            |
| 1 3 3 0 | ファイルシステム・サーバ        |
| 1 3 4 0 | ストレージ・エリア・ネットワーク    |
| 1 3 5 0 | ストレージ・プール           |
| 1 6 0 0 | 読取マッピング・テーブル        |
| 1 6 1 0 | 書込マッピング・テーブル        |

10

20

## 【图 1】



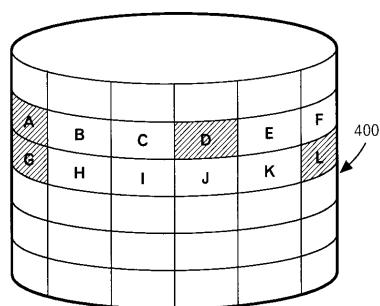
## 【圖 2】

範囲／ オフセット	仮想／相対 ブロック番号
0－4K	1
4－8K	2
8－12K	3
12－16K	4

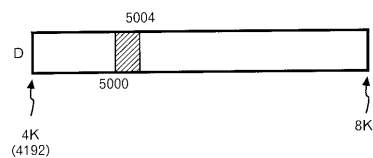
## 【 圖 3 】

300 ファイル・システム ・マッピング・テーブル	
仮想／相対 ブロック番号	物理ブロック ・アドレス
1	A
2	D
3	G
4	L

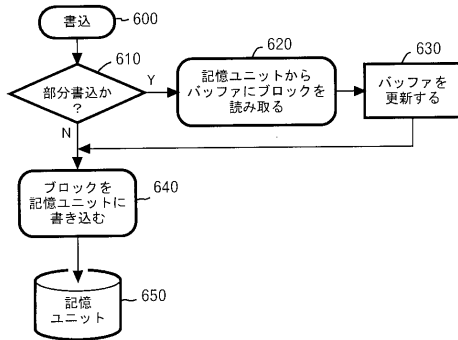
## 【图 4】



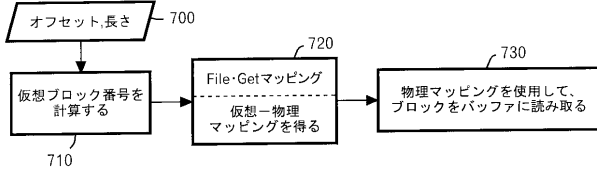
## 【 図 5 】



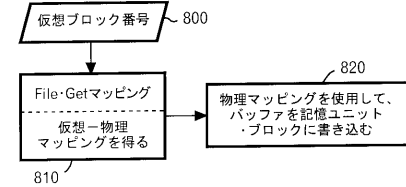
【図 6】



【図 7】



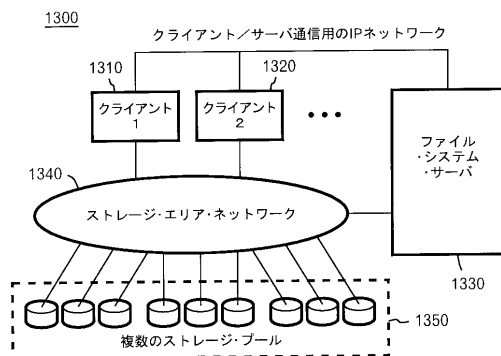
【図 8】



【図 12】

1200 読取 マッピング・テーブル		1210 書込 マッピング・テーブル	
仮想 ブロック 番号	物理ブロック 番号	仮想 ブロック 番号	物理ブロック 番号
1	A	1	W
2	X	2	X
3	G	3	Y
4	L	4	Z

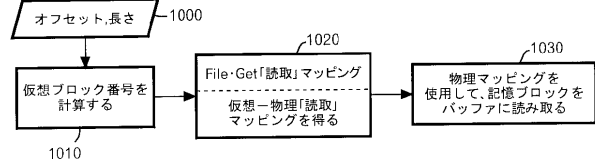
【図 13】



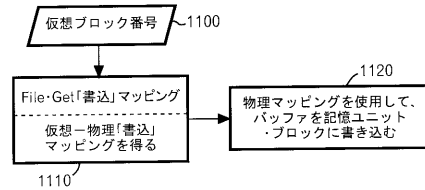
【図 9】

900 読取 マッピング・テーブル		910 書込 マッピング・テーブル	
仮想 ブロック 番号	物理ブロック 番号	仮想 ブロック 番号	物理ブロック 番号
1	A	1	W
2	D	2	X
3	G	3	Y
4	L	4	Z

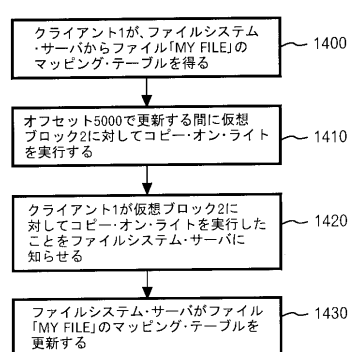
【図 10】



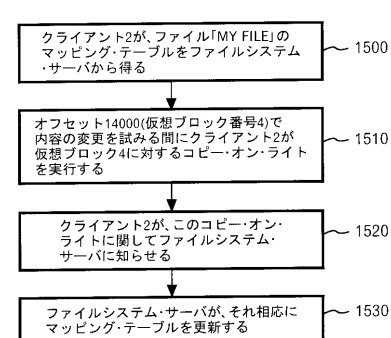
【図 11】



【図 14】



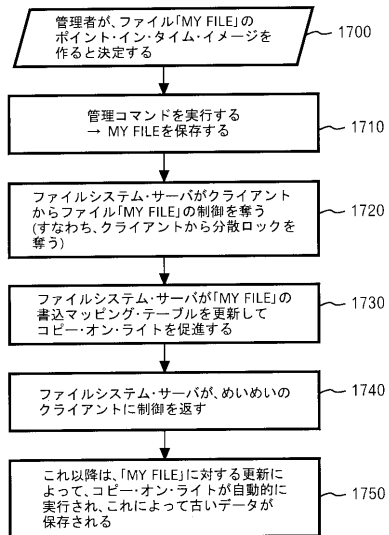
【図 15】



【図 16】

1600 読取 マッピング・テーブル		1610 書込 マッピング・テーブル	
仮想 ブロック番号	物理ブロック 番号	仮想 ブロック番号	物理ブロック 番号
1	A	1	W
2	X	2	X
3	G	3	Y
4	Z	4	Z

【図 17】



## フロントページの続き

- (72)発明者 ラジャゴパル・アナンタナラヤナン  
アメリカ合衆国 9 5 0 3 5 カリフォルニア州ミルピータス サーク・コート 5 2 2
- (72)発明者 ラルフ・エー・ベッカーツェンディ  
アメリカ合衆国 9 5 0 3 3 カリフォルニア州ロス・ガトス サンセット・リッジ・ロード 7 3  
5
- (72)発明者 ロバート・エム・リース  
アメリカ合衆国 9 5 0 3 3 カリフォルニア州ロス・ガトス グリーンウッド・ドライブ 1 7 9  
6 3
- (72)発明者 ランダル・シー・バーンズ  
アメリカ合衆国 2 0 0 1 1 ワシントン D . C . ウェブスター・ストリート ノース・ウェスト  
1 5 2 0
- (72)発明者 ダレル・ディー・イー・ロング  
アメリカ合衆国 9 5 0 7 3 カリフォルニア州ソーケル ワイルダー・ドライブ 5 0 3 0
- (72)発明者 ヴェンカテスワラオ・ジュジュリ  
アメリカ合衆国 9 7 0 0 6 オレゴン州ビーヴァートン ノースウェスト・トゥーソン・ストリー  
ト 1 7 1 7 7
- (72)発明者 デーヴィッド・エム・ウルフ  
アメリカ合衆国 9 7 2 2 9 オレゴン州ポートランド ノースウェスト・スターク・コート 9 1  
9 5
- (72)発明者 ジェイソン・シー・ヤング  
アメリカ合衆国 9 7 2 2 9 オレゴン州ポートランド ノースウェスト・カントリリッジ・ドライ  
ブ 1 6 9 7 2
- F ターム(参考) 5B065 BA01 CC01 CC03 EA33  
5B082 FA05